



The 4th International Conference on Arabic Computational Linguistics (ACLing 2018),  
November 17-19 2018, Dubai, United Arab Emirates

## Towards an Optimal Solution to Lemmatization in Arabic

Abed Alhakim Fraihat<sup>a,\*</sup>, Mourad Abbas<sup>b</sup>, Gábor Bella<sup>a</sup>, Fausto Giunchiglia<sup>a</sup>

<sup>a</sup>Department of Information Engineering and Computer Science, University of Trento, via Sommarive, 5, 38123 Trento, Italy

<sup>b</sup>Computational Linguistics Department, CRSTDLA, Algeria

### Abstract

Lemmatization—computing the canonical forms of words in running text—is an important component in any NLP system and a key preprocessing step for most applications that rely on natural language understanding. In the case of Arabic, lemmatization is a complex task because of the rich morphology, agglutinative aspects, and lexical ambiguity due to the absence of short vowels in writing. In this paper, we introduce a new lemmatizer tool that combines a machine-learning-based approach with a lemmatization dictionary, the latter providing increased accuracy, robustness, and flexibility to the former. Our evaluations yield a performance of over 98% for the entire lemmatization pipeline. The lemmatizer tools are freely downloadable for private and research purposes.

© 2018 The Authors. Published by Elsevier B.V.

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/3.0/>)

Peer-review under responsibility of the scientific committee of the 4th International Conference on Arabic Computational Linguistics.

**Keywords:** lemmatization; Arabic; natural language processing; machine-learning-based lemmatization; dictionary-based lemmatization

### 1. Introduction

Lemmatization consists of assigning to the surface form of each word in a text its corresponding *lemma*, that is, its canonical form as the word is commonly found in a dictionary. As such, lemmatization decreases morphological variations in text, in turn facilitating operations such as semantic analysis [1], information retrieval [2], question answering [3], or search [4]. For this reason, lemmatization is a crucial preprocessing operation in a wide range of applications that involve dealing with natural language.

The difficulty of the lemmatization task greatly depends on the nature of the language. In morphologically poor languages such as English, lemmatization can be considered an easy task already solved by simple normalization rules and a list of exceptional cases. For morphologically rich—highly inflecting or agglutinative—languages such as Arabic, on the other hand, it remains difficult and requires diverse, more complex approaches that are often specific to the language [5].

The major challenges specific to Arabic lemmatization, and NLP in general, are the rich morphology [6], which includes agglutinative properties [7], and the optional—and mostly omitted—marking of short vowels in writing [8]. This last property results in pervasive lexical ambiguity even considering the corresponding part of speech: for example, the past tense verb *صرت* could be vocalized as *صِرْتُ* with the corresponding lemma *صَارَ/become* or as *صَرَّ* with

\* Corresponding author. Tel.: +0-000-000-0000 ; fax: +0-000-000-0000.

E-mail address: [abdel.fraihat@gmail.com](mailto:abdel.fraihat@gmail.com)

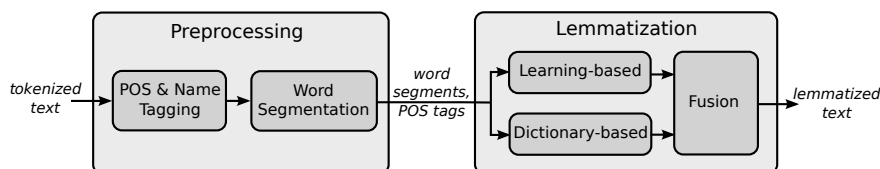


Fig. 1. The high-level NLP pipeline architecture for lemmatization

the corresponding lemma *صَرَ/ɡrit*. As deciding for the correct lemma is ultimately a word sense disambiguation problem, such cases put considerable stress on the quality of lemmatization. Tools that are capable of outputting multiple solutions in an order of preference are in this sense more robust as they potentially allow the disambiguation problem to be delayed to later, syntactic or semantic processing steps.

There has been extensive research so far on solving the lemmatization problem for Arabic. While several approaches were proposed, there are no more than a handful of actual tools available. Existing tools typically combine multiple techniques to achieve efficient lemmatization. The *Alkhalil* lemmatizer [9] first applies morpho-syntactic analysis to the input sentence to generate all potential word surface forms. Then, among these only one form is selected per word using a technique based on hidden Markov models. The accuracy of the tool is reported to be about 94%. Another lemmatizer is *Madamira* [10] which relies on preliminary morphological analysis on the input word that outputs a list of possible analyses. As a second step, it predicts the correct lemma using language models. The accuracy of the tool is 96.6%. The *Farasa* lemmatizer [11] uses a dictionary of words and their diacritizations ordered according to their number of occurrences. The accuracy reported for Farasa is 97.32%.

Beside these tools, there are other proposed approaches: for example, [12] propose a pattern-based approach while [13] and [14] present rule-based solutions.

In this paper, we present a new, freely available lemmatization tool that is composed of the fusion of a machine-learning-based classifier as a main lemmatizer and of an auxiliary dictionary-based lemmatizer. The underlying idea is that the classifier, that relies on context, is well-suited to solving cases of lexical ambiguity while the dictionary-based extension provides an extra performance boost, easy extensibility by new lemmas (e.g., neologisms), as well as the possibility to retrieve multiple possible lemmas per word form for subsequent analysis. The two lemmatizers are, however, implemented as separate tools and can also be used independently, each one yielding results beyond 95% of accuracy. Their high performance is partly explained by the preceding, fast and lightweight steps of POS tagging, morphological analysis, and word segmentation—earlier contributions of the authors reused in this work—that provide rich morphological information to the lemmatizers.

The lemmatizer was implemented as a component in a new, free, comprehensive pipeline for Arabic NLP [15] and is freely available for private or research purposes.<sup>1</sup> Both the 3-million-entry dictionary and the 2-million-token annotated corpus used to train the classifier were entirely generated by the authors and are contributions of this paper.

## 2. The Lemmatization Pipeline

As is the case of most approaches, our lemmatizer operates over an input pre-annotated by previous preprocessing steps. The pipeline specific to our method is shown in figure 1 and is composed of the following main steps:

1. *Preprocessing*: taking whitespace-tokenized Arabic text in input, we pre-annotate the text through the following operations:
  - (a) *POS and name tagging*: tokens are annotated by a machine-learning-based sequence labeler that outputs both POS and named entity tags, later used by the lemmatizer;
  - (b) *word segmentation*: using the POS output, cliticized words are segmented into a proclitic, a base word, and an enclitic, making the subsequent lemmatization step simpler.

<sup>1</sup> <http://www.arabicnlp.pro/alp/>

2. *Lemmatization*: the segmented and pre-annotated text is fed into the following lemmatizer components:

- (a) *dictionary-based lemmatizer*: words are lemmatized through dictionary lookup;
- (b) *machine-learning-based lemmatizer*: words are lemmatized by a trained machine learning lemmatizer;
- (c) *fusion*: the outputs of the two lemmatizers are combined into a single output.

In the following sections we present each component in detail, with a focus on lemmatization components as the pre-processors have already been discussed in our earlier work [15].

### 3. Preprocessing

The role of preprocessing is to enrich the input of the lemmatizer (and other subsequent components) by morphological and other contextual information, based on which the lemmatization task is simplified. With respect to state-of-the-art tools [9, 10], the preprocessing required by our lemmatization approach is lightweight and fast. The tokenized input text is first enriched by part-of-speech tagging, implemented as a fast machine-learning-based sequence labeler. Then it is further segmented by a very simple word segmenter component that further reduces the complexity and ambiguity of words. The simplicity of the process, presented in detail in [15], is explained by the rich morphological information output by the POS tagger, based on which word segmentation becomes a nearly trivial task. For example, the word *وباستخدام* (*and by using*) is first identified by the POS tagger to contain two proclitics and an inflected base word. On the basis of this result the segmenter outputs the word segment sequence <ب, و, استخدام> (<*and, by, using*>) and the corresponding POS tags <C, P, SMN> (<conjunction, preposition, singular masculine noun>).

In this section we provide a brief overview of the preprocessing from the point of view of the subsequent lemmatization task.

#### 3.1. POS and Name Tagging

In the lemmatization pipeline, the main goal of the POS and name tagger is to reduce the ambiguity of words by extracting information from their morphology and context:

- whether the word is part of a name or not;
- the corresponding part of speech;
- whether the word contains proclitics or enclitics (prefixes or suffixes).

The tagger is implemented as a single machine learning component, described in our earlier work [15] and freely available online.<sup>2</sup> In the following we provide only a brief presentation of it, in order to demonstrate the level of detail it provides to the subsequent lemmatizer.

```
<TAG>      ::= <PREFIX> <BASETAG> <POSTFIX>
<BASETAG> ::= <POSTAG> | <NERTAG>
<PREFIX>  ::= <PREFIX> | <PROCLITIC> "+" | ""
<POSTFIX> ::= <POSTFIX> | "+" <ENCLITIC> | ""
```

A tag is thus composed of a mandatory base tag and of zero or more proclitics and enclitics concatenated with the “+” sign indicating word segments. A base tag, in turn, is either a POS tag or a named entity (NER) tag. We do not consider name tags in the rest of the paper as they are irrelevant for lemmatization beyond the fact that names are skipped by the lemmatizer.

On a coarse-grained level, POS tags are divided into the following categories:

<sup>2</sup> <http://www.arabicnlp.pro/alp/>



يوم	SMN	يوم	مواطنین	PMN	مواطن
يوما	SMN	يوم	و	C	و
يومان	DMN	يوم	إقلاق	SMN	إقلاق
يوما	DMN	يوم	ال	D	ال
يومين	DMN	يوم	راحة	SFN	راحة
يومي	DMN	يوم	ال	D	ال
أيام	PTN	يوم	عامّة	SFAJ	عام
أياماً	PTN	يوم	في	P	في

(a)

(b)

Fig. 4. Examples of the contents of (a) the dictionary and (b) of the corpus used to train the classifier.

approach is justified by the inherent ambiguity of diacritic-free Arabic words whose meanings are typically deduced, by humans and machines alike, from context. While the preliminary POS tagging resolves a great deal of ambiguity, some cases still remain such as the verb form *يكون* which may be the verb form of the verb *كان*, or the verb *كون*.

A downside of learning-based lemmatization is that more rare and exceptional cases, such as *أسنة* (*spears*), may not be covered by its training corpus, which leads to lemmatization mistakes. The addition of new cases requires the re-training of the classifier. Another inconvenience is that classifiers—such as OpenNLP that we used—typically commit on a single output result, which may or may not be correct. In case of such ambiguity, from the full set of possible lemmas further NLP processing steps may be able to provide a correct results based on, e.g., syntactic or semantic analysis. In order to support these cases, we complement the learning-based lemmatizer by a dictionary-based one. The dictionary lemmatizer can be run independently, but we also provide a simple fusion method that combines the results of the two lemmatizers as described below.

Both lemmatizer components were implemented using the Apache machine learning based toolkit OpenNLP,<sup>4</sup> using the Maximum Entropy classifier.

#### 4.1. Dictionary Lemmatization

The dictionary of a db lemmatizer consists of a text file containing, for each row, a word, its POS tag and the corresponding lemma, each column separated by a tab character. An example of the of the dictionary is shown in column (A) of figure 4.

In case of ambiguous word forms (i.e., a word form POS-tag pair that has several lemmas), the corresponding lemmas are separated by “#” character. For example lemmas of the word form *زور* are *زار#زور*.

In the following we describe the method we used to build the dictionary. The used corpus is the same corpus we used for segmenting, POS-tagging, and named entity recognition as described in our previous work [15].

1. **Segmentation:** The corpus was segmented as explained in the previous section. The result of this step was generating a segmented corpus that contains more than 3.1 million segmented tokens.
2. **POS-tag based classification:** In this step, we classified the word forms according to their POS-tag.
3. **Inherent feminine and adjectival feminine classification:** In this step, we classified the feminine nouns into inherent feminine and adjectival feminine nouns. For example, the noun *أسرة*.SFN/*family* is inherent feminine while the noun *أسيرة*.SFN/*prisoner* is adjectival. This differentiation is important because the lemma of adjectival nouns is the masculine singular form of the noun *أسير* while there is no masculine singular lemma for *أسيرة*.
4. **Plural type classification:** In this step, we classified the singular, and dual nouns (after extracting their singular forms) according to their plural type into six classes as shown in Table 1. This classification enables us to build the possible word number\_gender forms of a given lemma automatically. For example, the class SMN\_PMN has six different possible number\_gender forms. On the other hand, using the feminine classification lists in previous step, enabled us to differentiate between the SMN\_PFN and SFN\_PFN. In the class SMN\_PFN, the lemma of a singular feminine noun (SFN) is the singular masculine noun (SMN). In the class SFN\_PFN on the other hand,

<sup>4</sup> <http://opennlp.apache.org/docs/1.9.0/manual/opennlp.html#tools.cli.lemmatizer>

the lemma is the singular feminine noun itself. The adjectives were classified into three classes. The first class is similar to the class SMN\_PMN which allows six different word forms. The second class contains a seventh possible form which is the broken plural adjective form. The third class contains PIAJ as a single possible plural form. For example, *ماهر* belongs to this class since it has two possible plural forms *ماهرون* and *مهرة*.

5. **Lemmas extraction:** This step is semi-automatic as follows:

- **Manual:** Assigning the lemmas to broken plural nouns and adjectives was performed manually.
- **Automatic:** Based on the morphological features in the tags, it was possible to extract lemmas for singular, dual, masculine plurals, feminine plurals adjectives and nouns. We also used rules to extract the verb lemmas such as removing the affixes *وا*.

6. **lemmas Enrichment:** Using the lemmas from the previous step, we have enriched the corpus with new verbs, adjectives, and nouns. For example, if the lemma of a plural noun or adjective was missing, we added it to the noun and adjective lemmas lists.

7. **Dictionary generation:** The files produced so far are as follows.

- Noun files:** Three files for masculine, feminine, and foreign nouns. The lemmas in these files were classified according to Table 1. There is a fourth file that contains quantifiers, pronouns, adverbs, ...
- Adjectives:** Three files for adjectives, comparatives, and ordinal adjectives. The lemmas in the adjective file are classified according to Table 1.
- Verbs:** One file that contains all extracted verb lemmas.

Using these files, the dictionary was generated as follows:

- **Nouns and adjectives generation:** According to the plural class, the noun and adjective forms were generated. The *ا* case ending, or changing *ة* to *ت* were also considered in this step.
- **Verbs generation:** For each verb in the verb lemmas list, we automatically generated the verb conjugations in present, past, imperative cases. We considered also accusative (*فعل منصوب*) and asserted verbs (*فعل مجزوم*).
- **Dictionary building:** Using the results from previous step, we built the dictionary as described in Figure 4 (B), where the lemmas of ambiguous surface forms were joined into a single string using the # operator.

Table 1. Plural classes

Class	Possible Word Forms	Example
SMN_PMN	SMN,SFN,DMN, DFN, PFN, PMN	مؤمنون, مؤمنات, مؤمنتان, مؤمنان, مؤمنة, مؤمن, مؤمن
SMN_PFN	SMN,DMN,PFN	إجراءات, إجراءات, إجراء, إجراء, إجراء
SFN_PFN	SFN,DFN,PFN	كتابات, كتابتين, كتابة, كتابة
SMN_PIN	SMN,DMN,PIN	أعلام, علمان, علم, علم
SFN_PIN	SFN,DFN,PIN	أسر, أسر, أسر, أسر
FWN_PFN	FWN,DMN,PFN	تلفزيونات, تلفزونان, تلفزيون, تلفزيون

#### 4.2. Machine Learning Lemmatization

The format of the corpus here is similar to the dictionary described in previous section. The only difference is that we order the entries according to their original position in the sentence in the segmented corpus. An empty line

indicates the end of a sentence. An example of the of the training corpus is shown in column (B) of figure 4. We used the segmented corpus from previous section to build the lemmatization corpus was performed in two steps:

1. **Lemmas assignation:** In this step, we used a dictionary lemmatizer to assign the word forms to their corresponding lemmas. In case of preposition, particles, and numbers, the lemma of the word form was a normalized form of the word form itself. The lemmas of named entities were also the named entities themselves. In this step, if a word form was ambiguous, all its possible lemmas were assigned.
2. **Validation:** In this step, we disambiguated the lemmas of the ambiguous word forms Manually.

The size of the generated corpus is 3,229,403 lines. The unique word forms after discarding the digits is 59,049 as specified in Table 2.

Table 2. Lemmas and unique word forms distribution in the corpus of the mlb lemmatizer

POS	Number of lemmas	Number of word forms
Noun	18,165	26,337
Adjective	6,369	13,703
Verb	4,258	19,009
Named entity	20,407	20,407
Particle	605	649

In a final step, we added all generated word forms and their corresponding lemmas from the dictionary described in previous section to the corpus. This increased the size of the corpus to 3,890,737 lines.

### 4.3. Fusion

While the learning-based lemmatizer outputs for each word a single candidate lemma, from the dictionary multiple solutions could be retrieved even for a single part of speech (for example, the verb form *يسلم* can be a verb form of the verb *سلم* gave up, or *أسلم* converted to Islam). The goal of the simple fusion component is to produce a final result from these solutions. The final output is a list of one or more lemmas in a decreasing order of confidence.

The idea underlying the fusion method is that we *usually* trust the dictionary to be capable of providing a correct solution space (a small set of possible lemmas), while we *usually* trust the classifier to return the most likely lemma from the previous set. However, in the case of out-of-corpus words the classifier may return incorrect results extrapolated from similar examples, such as returning the lemma *تہمہن* for the word form *یتہمہن*. Thus, whenever a lemma is returned by the classifier that is not included in the dictionary, it will still be included as a solution but with a lower confidence.

Accordingly, our simple fusion method is as follows. We take as input the results output by the two lemmatizers, namely  $L_{DIC} = \{l_1, \dots, l_n\}$  for the dictionary-based one and  $L_{CL} = \{l\}$  for the classifier-based one, and output  $L_F$ , the fusion result. We start by comparing the results of the two lemmatizers:

- if  $|L_{DIC}| = 1$  and  $l_1 = l$ , i.e., the outputs are identical, then the solution is trivial, we return either output and we are done:  $L_F = \{l\}$ ;
- otherwise, two further cases are distinguished:
  - if  $l \in L_{DIC}$ , that is, the dictionary contains the classification output, then we prioritize the result of the classifier by making it first (i.e., the preferred lemma):  $L_F = \{l, l_1, \dots, l_n\}$ ;
  - otherwise, we add the classifier result as the last element:  $L_F = \{l_1, \dots, l_n, l\}$ .

## 5. Evaluations

For evaluation we used a corpus of a 46,018-token text, retrieved and assembled from several news portals (such as Aljazeera news portal<sup>5</sup> and Al-quds Al-Arabi news paper<sup>6</sup>). We excluded from the evaluation the categories of tokens that cannot be lemmatized: 5,853 punctuation tokens, 3,829 tokens tagged as named entities, 482 digit tokens, and 10 malformed tokens (i.e., containing typos, such as *أوروبية بالموافق* instead of *أوروبية بالموافق*). Thus the number of tokens considered was 35,844.

In order to have a clear idea of the efficiency of the lemmatization pipeline, we evaluated it in a fine-grained manner, manually classifying the mistakes according to the component involved. This allows us to compute a comprehensive accuracy for the entire pipeline as well as evaluate individual components: the POS tagger, the segmenter, each lemmatizer, as well as the fusion lemmatization. The evaluation data files are available online.<sup>7</sup>

Table 3. Types of mistakes committed by the learning-based lemmatizer, and their proportions

Type of mistake	Occurrences	Example
POS tag (coarse-grained) mistakes	199	تبريح.SMN instead of تبريح.PRSV
Morphological tag (fine-grained) mistakes	201	كروم.SMN instead of كروم.PIN
Segmentation tag mistakes	103	أخذتنا.PSTV instead of أخذت وأخذنا.PRO
Classifier mistakes: nonexistent lemma	158	جد instead of وجد for يجدوا.PRSV
Classifier mistakes: wrong disambiguation	12	كان instead of كون for يكونوا.PRSV
Dictionary mistakes: missing word form	1,207	كاشطة, دواعم, قرفضاء
Fusion mistakes	50	نلن, عوالم, دواعم

The fine-grained evaluation is summed up in table 3.<sup>8</sup> *Nonexistent lemma* stands for cases where the POS tag and the segmentation were correct, yet the classifier gave a wrong, non-linguistic result. *Wrong disambiguation* means that the lemmatizer chose an existing but incorrect lemma for an ambiguous word form.

Table 4. Accuracy values computed for various components of the lemmatization pipeline

Component	Evaluation method	Accuracy
preprocessing	all mistakes (POS, morphological, segmentation)	98.6%
classifier-based lemmatizer	in isolation	99.5%
classifier-based lemmatizer	in isolation, built-in OpenNLP cross-validation	99.7%
classifier-based lemmatizer	entire pipeline	98.1%
dictionary-based lemmatizer	in isolation	96.6%
dictionary-based lemmatizer	entire pipeline	95.2%
fusion lemmatizer	entire pipeline	98.4%

The accuracy measures reported in table 4 were computed based on the results in table 3. On these we make the following observations. The performance of preprocessing (98.6%) represents an upper bound for the entire lemmatization pipeline. In this perspective, the near-perfect results of the classifier (99.5% when evaluated in isolation, 98.1% on the entire pipeline) are remarkable. We cross-checked these results using the built-in cross-validation feature of

<sup>5</sup> <http://www.aljazeera.net/>

<sup>6</sup> <http://www.alquds.co.uk/>

<sup>7</sup> <http://www.arabicnlp.pro/alp/lemmatizationEval.zip>

<sup>8</sup> While after tagging and segmentation the number of (segmented) tokens rose to 62,694, we computed our evaluation results based on the number of unsegmented tokens.



OpenNLP and obtained similar results (99.7%). The dictionary-based lemmatizer reached a somewhat lower yet still very decent result (96.6% in isolation, 95.2% on the entire pipeline), due to the 1207 OOV word forms. The fusion of the two lemmatizers, finally, improved slightly on the classifier: of the 170 mistakes made by the classifier, 120 could be correctly lemmatized using the dictionary. Thus the fusion method reached a full-pipeline result of 98.4%, only a tiny bit worse than the performance of preprocessing itself.

## 6. Conclusion and Future Work

We presented an optimization approach for Arabic lemmatization, based on the combination of machine learning and a lemmatization dictionary, that provides excellent accuracy. Beside the result itself, the addition of a lemmatization dictionary provides additional robustness to the underlying NLP pipeline. Firstly, it makes the lemmatizer easy to extend by new lemmas that could potentially be mislabeled by the classifier. Secondly, it allows the lemmatizer to return not only one result but an order list of candidate lemmas, allowing the decision to be delayed to subsequent NLP components.

Both the machine learning model and the dictionary were built using a corpus of 2.2 million tokens annotated and manually validated by the authors. The dictionary, the trained model, and corresponding tools are all free for research purposes upon request.

The presented tool was implemented as a component of the *ALP* comprehensive NLP pipeline. We plan to extend the current pipeline with new components such as a vocalizer, a phrase chunker, a dependency parser, or a multiword expression detector.

## References

- [1] R. Navigli, "Word sense disambiguation: A survey," *ACM Comput. Surv.*, vol. 41, pp. 10:1–10:69, Feb. 2009.
- [2] V. Balakrishnan and L.-Y. Ethel, "Stemming and lemmatization: A comparison of retrieval performances," vol. 2, pp. 262–267, 01 2014.
- [3] A. A. Freihat, M. R. H. Qwaider, and F. Giunchiglia, "Using grice maxims in ranking community question answers," in *Proceedings of the Tenth International Conference on Information, Process, and Knowledge Management, eKNOW 2018, Rome, Italy, March 25-29, 2018*, pp. 38–43, 2018.
- [4] F. Giunchiglia, U. Kharkevich, and I. Zaihrayeu, "Concept search," in *The Semantic Web: Research and Applications* (L. Aroyo, P. Traverso, F. Ciravegna, P. Cimiano, T. Heath, E. Hyvönen, R. Mizoguchi, E. Oren, M. Sabou, and E. Simperl, eds.), (Berlin, Heidelberg), pp. 429–444, Springer Berlin Heidelberg, 2009.
- [5] A. Farghaly and K. Shaalan, "Arabic natural language processing: Challenges and solutions," vol. 8, pp. 14:1–14:22, Dec. 2009.
- [6] M. Gridach and N. Chenfour, "Developing a new approach for arabic morphological analysis and generation," *CoRR*, vol. abs/1101.5494, 2011.
- [7] K. Shaalan, "A survey of arabic named entity recognition and classification," *Comput. Linguist.*, vol. 40, pp. 469–510, June 2014.
- [8] O. Hamed and T. Zesch, "A Survey and Comparative Study of Arabic Diacritization Tools," *JLCL: Special Issue - NLP for Perso-Arabic Alphabets.*, vol. 32, no. 1, pp. 27–47, 2017.
- [9] M. Boudchiche, A. Mazroui, M. Ould Abdallahi Ould Bebah, A. Lakhouaja, and A. Boudlal, "Alkhalil morpho sys 2: A robust arabic morpho-syntactic analyzer," *J. King Saud Univ. Comput. Inf. Sci.*, vol. 29, pp. 141–146, Apr. 2017.
- [10] A. Pasha, M. Al-Badrashiny, M. T. Diab, A. El Kholly, R. Eskander, N. Habash, M. Pooleery, O. Rambow, and R. Roth, "Madamira: A fast, comprehensive tool for morphological analysis and disambiguation of arabic.," in *LREC*, vol. 14, pp. 1094–1101, 2014.
- [11] A. Abdelali, K. Darwish, N. Durrani, and H. Mubarak, "Farasa: A fast and furious segmenter for arabic," in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pp. 11–16, Association for Computational Linguistics, San Diego, California, 2016.
- [12] M. Attia, A. Zirikly, and M. T. Diab, "The power of language music: Arabic lemmatization through patterns," in *Proceedings of the 5th Workshop on Cognitive Aspects of the Lexicon, CogALex@COLING 2016, Osaka, Japan, December 12, 2016*, pp. 40–50, 2016.
- [13] E. Al-Shammari and J. Lin, "A novel arabic lemmatization algorithm," in *Proceedings of the Second Workshop on Analytics for Noisy Unstructured Text Data, AND '08, (New York, NY, USA)*, pp. 113–118, ACM, 2008.
- [14] T. El-Shishtawy and F. El-Ghannam, "An accurate arabic root-based lemmatizer for information retrieval purposes," *CoRR*, vol. abs/1203.3584, 2012.
- [15] A. A. Freihat, G. Bella, H. Mubarak, and F. Giunchiglia, "A single-model approach for arabic segmentation, pos tagging, and named entity recognition," in *2018 2nd International Conference on Natural Language and Speech Processing (ICNLSP)*, pp. 1–8, April 2018.
- [16] I. Zeroual, A. Lakhouaja, and R. Belahbib, "Towards a standard part of speech tagset for the arabic language," *Journal of King Saud University - Computer and Information Sciences*, vol. 29, no. 2, pp. 171 – 178, 2017. Arabic Natural Language Processing: Models, Systems and Applications.