

ARTICLE

UNLT: Urdu Natural Language Toolkit

Jawad Shafi^{1,2,*}, Hafiz Rizwan Iqbal³, Rao Muhammad Adeel Nawab² and Paul Rayson¹

¹School of Computer and Communication (SCC), Lancaster University, Lancaster, UK, ²COMSATS University Islamabad, Lahore Campus, Pakistan and ³Information Technology University, Lahore, Pakistan

*Corresponding author. E-mail: jawadshafi@cuilahore.edu.pk

(Received 31 December 2017; revised 15 November 2021; accepted 22 November 2021;
first published online 19 January 2022)

Abstract

This study describes a Natural Language Processing (NLP) toolkit, as the first contribution of a larger project, for an under-resourced language—Urdu. In previous studies, standard NLP toolkits have been developed for English and many other languages. There is also a dire need for standard text processing tools and methods for Urdu, despite it being widely spoken in different parts of the world with a large amount of digital text being readily available. This study presents the first version of the *UNLT (Urdu Natural Language Toolkit)* which contains three key text processing tools required for an Urdu NLP pipeline; word tokenizer, sentence tokenizer, and part-of-speech (POS) tagger. The UNLT word tokenizer employs a morpheme matching algorithm coupled with a state-of-the-art stochastic n -gram language model with back-off and smoothing characteristics for the space omission problem. The space insertion problem for compound words is tackled using a dictionary look-up technique. The UNLT sentence tokenizer is a combination of various machine learning, rule-based, regular-expressions, and dictionary look-up techniques. Finally, the UNLT POS taggers are based on Hidden Markov Model and Maximum Entropy-based stochastic techniques. In addition, we have developed large gold standard training and testing data sets to improve and evaluate the performance of new techniques for Urdu word tokenization, sentence tokenization, and POS tagging. For comparison purposes, we have compared the proposed approaches with several methods. Our proposed UNLT, the training and testing data sets, and supporting resources are all free and publicly available for academic use.

Keywords: Word segmentation; Text segmentation; Part-of-speech tagging; Urdu NLP toolkit; Urdu corpora

1. Introduction

A natural language toolkit is a library or framework used to analyze human language in a statistical, rule-based or hybrid Natural Language Processing (NLP). These toolkits have been used in the development of a range of applications from various domains. For instance, life sciences and medicine (Cunningham *et al.* 2013), bio-informatics (Ferrucci and Lally 2004), computer science (Maynard *et al.* 2015; Rush, Chopra, and Weston 2015), linguistics (Gries and John 2014), Machine Learning (ML) (Bird *et al.* 2008), and the analysis of social media (Dietzel and Maynard 2015).

The majority of existing NLP toolkits are for English with many other languages supported (Cunningham *et al.* 2002; Bird, Klein, and Loper 2009; Manning *et al.* 2014; Kwartler 2017). However, there is a lack of standard text processing tools and methods for South Asian languages, particularly Urdu, which has 300 million speakers around the world (Riaz 2009) for which a large amount of digital text is available through online repositories. Urdu is an Indo-Aryan^a (or Indic)

^ahttps://en.wikipedia.org/wiki/Indo-Aryan_languages#cite_note-ethnologue-4 - Last checked: 20-March-2019.



language derived from Sanskrit/Hindustani language (Bögel *et al.* 2007), has been heavily influenced by Arabic, Persian (Bögel *et al.* 2007) and less by Turkic (Chagatai^b) languages for literary and technical vocabulary (Sharjeel, Nawab, and Rayson 2017), and is written from right to left in Nastaliq style (Shafi 2020). Urdu is a morphologically rich language (Saeed *et al.* 2012), including many multiword expressions and letters which may change their shape based on context, which makes the tokenization task very complex and challenging. Moreover, it is a free word order language (Mukund, Srihari, and Peterson 2010; Riaz 2012; Daud, Khan, and Che 2016).

In the previous literature, a small amount of work has been carried out to propose systematic text processing NLP approaches for Urdu including word tokenization (Durrani and Hussain 2010; Lehal 2010; Rashid and Latif 2012; Rehman and Anwar 2012), sentence tokenization (Rehman and Anwar 2012; Raj *et al.* 2015), and POS tagging (Hardie 2004; Anwar *et al.* 2007a; b; Sajjad and Schmid 2009; Muaz, Ali, and Hussain 2009; Ahmed *et al.* 2014) (see Section 2 for details). However, there are a number of limitations of these existing Urdu NLP studies: (i) most of them are not formed into NLP tools and the ones that are implemented are not publicly and freely available, (ii) training/testing data sets along with developed resources are not always freely and publicly available to improve, compare, and evaluate new and existing methods (Daud *et al.* 2016), (iii) they have been trained and tested on very small data sets (Anwar *et al.* 2007a; Durrani and Hussain 2010; Rehman and Anwar 2012), (iv) the efficiency of dictionary-based word tokenization approach (Rashid and Latif 2012) is entirely dependent on a dictionary of complete Urdu words, which is not practically possible to produce for the Urdu language, (v) statistical *n*-gram based word tokenizers (Durrani and Hussain 2010; Rehman *et al.* 2013) cannot handle unknown words or back off to a lesser contextual models, (vi) sentence tokenization methods require data sets to train machine learning algorithms and these are unavailable (Rehman and Anwar 2012; Raj *et al.* 2015), (vii) current rule-based POS tagging methods (Hardie 2004) are closely tailored to a particular data set, therefore, not portable across different domains, (viii) smoothing and other features to handle unknown words in statistical POS taggers have not been thoroughly explored, (ix) less contextual POS tagging techniques have been proposed, and (x) POS tagsets that have been used to train/test statistical POS taggers (Hardie 2004; Anwar *et al.* 2007b; Sajjad and Schmid 2009) have several shortcomings (see Section 4.3.1).

To overcome the limitations of existing Urdu NLP studies, this study presents the UNLT (Urdu Natural Language Toolkit), initially with text processing methods and three NLP tools including a word tokenizer, sentence tokenizer, and POS tagger. Our proposed Urdu word tokenizer is based on a novel algorithm that makes use of rule-based morpheme matching, *n*-gram statistical model with backoff and smoothing characteristics, and dictionary look-up. It is trained on our proposed data set of 1361K tokens and evaluated on our proposed test data set containing 59K tokens. The UNLT sentence tokenizer is a combination of rule-based, regular expressions, and dictionary lookup techniques that are evaluated on our proposed test data set of 8K sentences. Furthermore, we have also proposed a ML-based sentence tokenizer. Finally, the UNLT POS tagger is mainly based on Hidden Markov and Maximum Entropy models, with multiple variations based on smoothing, suffix length, context window, word number, lexical, and morphological features. The set of UNLT POS taggers were trained on our proposed training data set of 180K tokens and evaluated on our proposed test data set of 20K tokens.

The UNLT and our training/testing data sets will be crucial in (i) fostering research in an under-resourced language, that is, Urdu, (ii) the development and evaluation of Urdu word tokenizers, sentence tokenizers, and POS taggers, (iii) facilitating comparative evaluations of existing and future methods for Urdu word tokenization, sentence tokenization, and POS tagging, (iv) the development of various NLP tools and applications in other areas such as information retrieval, corpus linguistics, plagiarism detection, semantic annotation, etc., and (v) providing a framework in which other Urdu NLP tools can be integrated.

^b<https://en.wikipedia.org/wiki/Urdu> - Last checked: 15-September-2020.

The rest of the article is organised as follows: Section 2 describes related work, and Section 3 presents challenges for word tokenization, sentence tokenization, and POS tagging methods. Section 4 explains the proposed UNLT modules. Whereas Section 5 present the proposed data sets. Section 6 introduces evaluation measures, results, and their analysis. Section 7 presents a summary and future directions of our research.

2. Related work

2.1 Existing Urdu word tokenization approaches

In the existing literature, we find only a few studies that have addressed the problem of word tokenization for the Urdu language, these are (Durrani and Hussain 2010; Lehal 2010; Rashid and Latif 2012; Rehman and Anwar 2012). The study in Rashid and Latif (2012) performs Urdu word tokenization in three phases. First, Urdu words are tokenized based on spaces, thus returning the cluster(s) of valid (single word) and invalid (merged word(s)^c) Urdu words. Next, a dictionary is checked against valid/invalid words to assure the robustness of the word(s). If the word is present in the dictionary, then it will be considered as a valid Urdu word, returning all single words. However, if the word is not matched in the dictionary, then it is considered as a merged word, hence, needing further segmentation. In the second phase, the merged words are divided into all possible combinations, to check the validity of each produced combination through dictionary lookup. If it is present in the dictionary, it will be considered as a valid word. The first two phases solve the problem of space omission (see Section 4.1), and the third phase addresses the space insertion problem by combining two consecutive words and checking them in the dictionary. If the compound word is found in the dictionary, then it will be considered as a single word. This technique of word tokenization was tested on 11,995 words with a reported error rate of 2.8%. However, the efficiency of this algorithm is totally dependent on the dictionary (used to check whether a word is valid or not), and it is practically not possible to have a complete dictionary of Urdu words. Furthermore, if a valid word is not present in the dictionary, then this technique will mark it as invalid, which will be wrong.

Durrani and Hussain (2010) proposed a hybrid Urdu word tokenizer that works in three phases. In the first phase, words are segmented based on space, thus returning a set of an orthographic word(s).^d Further, a rule-based maximum matching technique is used to generate all possible word segmentations of the orthographic words. In the second phase, the resulting words are ranked using minimum word heuristics, *uni*, and *bi*-grams based sequence probabilities. In the first two phases, the authors solved the space omission problem (see Section 4.1). In the third phase, the space insertion problem is solved to identify compound words by combining words using different algorithms. The proposed Urdu word tokenizer is trained on 70K words, whereas it is tested on a very small data set of 2367 words reporting an overall error rate of 4.2%. Although the authors have reported a very low error rate, this study has some serious limitations: (i) the evaluation is carried out on a very small data set, which makes the reported results less reliable in terms of how good the word tokenizer will perform on real-world data, (ii) using a statistical *n*-gram technique that may ultimately lead to data sparseness, and (iii) it does not tokenize Urdu text correctly even for short texts.

Another online CLE Urdu word tokenizer is available through a website,^e which allows tokenization of up to 100 words. Its implementation details are not provided. It reports an accuracy of 97.9%. However, the link is not always available,^f and its API is not freely available.^g We applied

^cCombination of many words.

^dOne orthographic word may eventually give multiple words and multiple orthographic words may combine to give a single word.

^e<http://182.180.102.251:8080/segment/> - Last checked: 24-June-2018.

^fSee: <http://182.180.102.251:8080/segment/> - As on: 11-April-2018.

^g<http://www.cle.org.pk/clestore/segmentation.htm> - Last checked: 24-June-2018.

the CLE online Urdu work tokenizer on three randomly selected input short texts, and they all were incorrectly tokenized with many mistakes.

The research cited in Lehal (2010) takes an approach to Urdu word tokenization, based on the Hindi language. The authors tokenized Urdu words after transliterating them from Hindi, as the Hindi language uses spaces consistently as compared to its counterpart Urdu. They also addressed and resolved the space omission problem for Urdu in two phases. In the first phase, Urdu grammar rules have been applied to decide if the Urdu adjacent words have to be merged or not. If the grammatical rules analyzer provides a definite answer that two adjacent words can be joined or not, then no further processing is required. However, if the rule-based analyzer is not confident about two words either it can be joined or not, then the second phase is invoked. In the second phase, Urdu and Hindi *uni*-gram and *bi*-gram bilingual lexical resources are used to make the final decision, that is, either we need to join the two adjacent words or not. This technique of Urdu word tokenization used 2.6 million words as training data, whereas it was tested on 1.8 million tokens. The results show an error rate of 1.44%. The limitations of this study are (i) the problem of space insertion has not been addressed, (ii) this approach requires large bilingual corpus which is difficult to create particularly for under-resourced languages like Urdu and Hindi.

Rehman *et al.* (2013) proposed an Urdu word tokenizer by using rule-based (maximum matching) with *n*-gram statistical approach. This approach to Urdu word tokenization uses several different algorithms to solve the problem of space omission and insertion. First, the forward maximum matching algorithm is used to return the list of individual tokens of Urdu text. Second, the Dynamic Maximum Matching (DMM) algorithm returns all the possible tokenized sequences of the Urdu text, segments are ranked and the best one is accepted. Third, DMM is combined with the *bi*-gram statistical language model. These three algorithms are used to solve the space omission problem, whereas for the space insertion problem, six different algorithms were used. The authors used 6400 tokens for training and 57,000 tokens for testing. This approach has produced up to 95.46% F_1 score. Furthermore, the algorithms are based on probabilities that may result in zero probability being assigned to some unknown words. The authors have not handled such cases with either back off or other smoothing estimators.

To overcome the limitations of the existing studies, our study proposes a novel Urdu word tokenization algorithm using a rule-based morpheme matching approach, with off-the-shelf statistical *tri*-gram language model with back-off and smoothing characteristics for the space omission problem, whereas the space insertion problem has been solved using dictionary lookup technique (see Section 4.1.2). To train and test our proposed word tokenizer, we developed benchmark training (contains 1.65 million tokens) and testing (contains 59K tokens) data sets. Our word tokenizer and training/testing data sets are freely and publicly available for download (see Section 7 for details).

2.2 Sentence tokenization approaches

The problem of Urdu sentence tokenization has not been thoroughly explored, and we found only two studies (Rehman and Anwar 2012; Raj *et al.* 2015) that address the issue. Rehman and Anwar (2012) used a hybrid approach that works in two stages. First, a *uni*-gram statistical model was trained on annotated data. The trained model was used to identify word boundaries on a test data set. In the second step, the authors used heuristic rules to identify sentence boundaries. This study achieved up to 99.48% precision, 86.35% recall, 92.45% F_1 , and 14% error rate, when trained on 3928 sentences; however, the authors did not mention any testing data. Although this study reports an acceptable score, it has some limitations: (1) the error rate is high (14%), (2) the evaluation is carried out on a very small data set, which makes the reported results less reliable, and it is difficult to tell how well the sentence tokenizer will perform on real test data, and (3) the trained model along with training/testing data are not publicly available.

In another study, Raj *et al.* (2015) used an Artificial Neural Network along with POS tags for sentence tokenization in two stages. In the first phase, a POS tagged data set is used to calculate the word-tag probability (P) based on the general likelihood ranking. Furthermore, the POS tagged data set along with probabilities was converted to bipolar descriptor arrays,^h to reduce the error as well as training time. In the next step, these arrays along with frequencies were then used to train feed forward Artificial Neural Network using back propagation algorithm and delta learning rules. The training and testing data used in this study are 2688 and 1600 sentences, respectively. The results show 90.15% precision, 97.29% recall, and 95.08% F_1 -measure with 0.1 threshold values. The limitations of this study are that the evaluation is carried out on a small set of test data, and the trained model, as well as the developed resources, is not publicly available.

Again, similar to the Urdu word tokenization problem (see Section 2.1), the developed Urdu sentence tokenizers along with training/testing data sets are not publicly available. To fill this gap, our contribution here is an Urdu sentence tokenizer which is a combination of rule-based, regular expressions and dictionary lookup techniques, along with training (contain 12K sentences) and testing (containing 8K sentences) data sets, all of which are free and publicly available for research purposes. Moreover, we have also proposed a novel supervised ML-based sentence tokenizer by extracting various features.

2.3 Part-Of-Speech tagging approaches

Similar to Urdu word (see Section 2.1) and sentence tokenization (see Section 2.2), the problem of Urdu POS tagging has not been thoroughly explored. We found only six studies (Hardie 2004; Anwar *et al.* 2007a; b; Sajjad and Schmid 2009; Muaz *et al.* 2009; Ahmed *et al.* 2014) that addressed the issue.

A pioneering piece of research on Urdu POS tagging is described in Hardie (2004). This work focused on the development of a uni-rule POS tagger, which consists of 270 manual crafted rules. The author used a POS tagset with 350 tags (Hardie 2003). The training data consist of 49K tokens, whereas testing was carried out on two different data sets containing 42K and 7K tokens. The reported average accuracy for the 42K tokens is 91.66%, whereas for the 7K corpus, the average accuracy is 89.26% with a very high ambiguity level (3.09 tags per word). However, the POS tagset used in this study has several limitations (see Section 4.3) and, therefore, cannot be used for a grammatical tagging task, and having a large number of POS tags with a relatively small training data will affect the accuracy and manually deducing rules is a laborious and expensive task.

The first stochastic POS tagger for the Urdu language was developed in 2007 (Anwar *et al.* 2007a). They have proposed a POS tagger based on a *bi*-gram Hidden Markov Model with back off to *uni*-gram model. Twoⁱ different POS tagsets were used. The reported average accuracies for the 250 POS tagset and 90 POS tagset are 88.82% and 92.60%, respectively. Both were trained on a data set of 1,000 words; however, the authors have not provided any information about the test data set. As before, this study has several limitations; the POS tagset of 250 tags has several grammatical deficiencies (see Section 4.3), the information about the proposed tagset of 90 tags is not available, the system was trained and tested on a very small data set, which shows that it is not feasible for morphological rich and free word order language, that is, Urdu, and used less contextual *bi/uni*-gram statistical models.

Anwar *et al.* (2007b) have developed an Urdu POS tagger using *bi*-gram Hidden Markov Model. The authors proposed six *bi*-gram Hidden Markov-based POS taggers with different smoothing techniques to resolve data sparseness. The accuracy of these six models varies from 90% to 96%. For each model, they used a POS tagset of 90 tags. However, the authors have not

^hIf $P > 0.1 \implies P \equiv -1$, If $P = 0.1 \implies P \equiv 0$, If $P < 0.1 \implies P \equiv +1$.

ⁱThe first POS tagset contains 250 POS tags (Hardie 2003), whereas the second one consists of 90 tags (details are not given).

mentioned the size of training/test data sets. This study has several limitations as, like the one in Anwar *et al.* (2007a), the authors have used a 350 POS tagset, which has several misclassifications (see Section 4.3), the training/testing data split is unknown to readers, and limited smoothing estimators have been used, it used *bi*-gram language model (i.e., less contextual), and suffix information has not been explored.

The authors in Sajjad and Schmid (2009) trained Trigrams-and-Tag (TnT) (Brants 2000), Tree Tagger (TT) (Schmid 1994b), Random Forest (RF) (Schmid and Laws 2008), and Support Vector Machine (SVM) (Giménez and Marquez 2004) POS taggers, using a tagset containing 42 POS tags. All these stochastic Urdu POS taggers were trained on a 100K word data set, whereas for testing only 9K words were used. The reported accuracy for TnT, TT, RF, and SVM is 93.40%, 93.02%, 93.28%, and 94.15%, respectively. In terms of limitations, they used a POS tagset of 42 tags, which has several grammatical irregularities (see Section 4.3).

In another study (Muaz *et al.* 2009), stochastic Urdu POS taggers are presented, that is, TnT and TT tagger. These taggers are trained and tested on two different data sets with the following statistics: (i) first data set consists of 101,428 tokens (4584 sentences) and 8670 tokens (404 sentences) for training and testing, respectively, and (ii) the second data set consists of 102,454 tokens (3509 sentences) and 21,181 tokens (755 sentences) for training and testing, respectively. The reported accuracy for the first data set is 93.01% for TnT tagger, whereas 93.37% for TT tagger. For the second data set, TnT tagger produced 88.13% accuracy and TT had 90.49% accuracy. Similar to other studies, it employed a POS tagset that has several grammatical problems (see Section 4.3), meaning that it is no longer practical for Urdu text.

The authors in Ahmed *et al.* (2014) have proposed an Urdu POS tagger^j which is based on Decision Trees and smoothing technique of Class Equivalence, using a tagset of 35 POS tags. It is trained and tested on the CLE Urdu Digest corpus,^k training and test data split are 80K and 20K tokens, respectively. However, this POS tagger is only available through an online interface, which allows tagging of 100 words. It is trained on a relatively small data set that is not freely available. The Decision Tree statistical models are less accurate for Urdu text as compared to HMM, etc. Sajjad and Schmid (2009) (see Section 6.3).

In contrast, our study contributes a set of Urdu POS taggers along with large training (containing 180K POS tagged tokens) and testing (containing 20K POS tagged tokens) data sets. Our proposed POS taggers are based on two machine learning techniques, *tri*-gram Hidden Markov Model (HMM) and Maximum Entropy (MaEn) based models. Each of our proposed HMM and MaEn models is a combination of different backoff, smoothing estimators, suffix, and other types of binary valued features. To the best of our knowledge, these models along with smoothing, backoff, suffix, and binary valued features have not been explored previously for the Urdu language.

3. Challenges of Urdu NLP tools

3.1 Challenges for word tokenization

Is a challenging and complex task for the Urdu language due to three main problems (Durrani and Hussain 2010): (i) the space omission problem—Urdu uses *Nastalique* writing style and *cursive script*, in which Urdu text does not often contain spaces between words, (ii) the space insertion problem—*irregular* use of spaces within two or more words, and (iii) ambiguity in defining Urdu words—in some cases, Urdu words lead to an ambiguity problem because there is no clear agreement to classify them as a single word or multiple words.

^j<http://182.180.102.251:8080/tag/> - Last checked: 09-July-2018.

^k<http://www.cle.org.pk/clestore/urdu Digestcorpus100k.htm> - Last checked: 09-July-2018.

Table 1. Example text for various types of space omission problems

Type	Correct	Incorrect	Translation
Affixation	خوش - اخلاق KHOSH AKHLAK	خوشاخلاق KHOSHAKHLAK	Polite
Abbreviations	این - ایل - ای AYN AYL AY	اینالای AYNAYLAY	NLE
Compound word	تغیر - پذیر TGHYR PZYR	تغیرپذیر TGHYR PZYR	Variable
English word	نٹ - ورک NYT ORK	نٹورک NYTORK	Network
Proper noun	ویسٹ - انڈیز OYST ANDYZ	ویسٹانڈیز OYSTANDYZ	West Indies
Reduplication	آنن - فائن AANN FANN	آننفاائن AANNFANN	Quickly

“new”), (ii) سائیکل (SAYYKL, “bicycle”), and (iii) ہے (HY, “is”). As it can be noted from the above discussion, space omission problems are complex thus making the Urdu word tokenization task particularly challenging.

In the space insertion problem, if the first word ends either on a joiner or nonjoiner, a space at the end of the first word (see Table 1, Correct column– incorrect multiple tokens with space (-), but correct shape) can be inserted for several reasons: (i) affixes can be separated from their root, (ii) to keep separate Urdu abbreviations when transliterated, (iii) increase readability for Urdu proper nouns and English/foreign words are transliterated, (iv) compound words and reduplication morphemes do not visually merge and form a correct shape, and (v) to avoid making words written incorrectly or from combining (see Table 1, incorrect column–single token but incorrect shape). For example, خوش اخلاق (KHOSH AKHLAK, “polite”) is a compound word of type affixation; however, space was inserted between خوش (KHOSH, “happy”), that is, a prefix (literally “happy”) and اخلاق (AKHLAK, “ethical”), that is, root to increase the readability and understandability. To identify خوش اخلاق (KHOSH AKHLAK, “polite”) as a single word/token the tokenizer will need to ignore the space between them. This also serves to emphasize the fact that the space insertion problem is also a very challenging and complex task in Urdu word tokenization.

As discussed earlier, in some cases, Urdu words are harder to disambiguate. There is no clear agreement on word boundaries in a few cases (sometimes they are considered as a single word even by a native speaker). For example the compound word, وزیر اعلیٰ (OZYR AALY, “chief minister”), بہن بھائی (BHN BHAYY, “sibling,” literally “brother sister”). The same is the case for reduplications, فر فر (FR FR, “fluent”) and affixation, بد اخلاق (BD AKHLAK, “depravedly”). Certain function words (normally case markers, postpositions, and auxiliaries) can be written jointly, for example, اس میں (AS MYN, “herein”), یہ وقت (YHOKT, “this time”), or ہوگی (Ho GEE). Alternatively, the same function words can be written separately such as اس میں (AS MYN, “herein”), یہ وقت (YH OKT, “this time”), and ہوگی (HO GEE) (i.e., two auxiliaries), respectively. These distinct forms of the same word(s) are visually correct and may be perceived as single or multiple words. These types of cases are ambiguous, that is, can be written with or without spaces

and can be treated as a single unit or two different words. Consequently, this changes the perception of where the word boundary should sit. A possible solution to handle such words is to use a knowledge base. To conclude, the space insertion problem, space omission problem, and ambiguity in tokenizing multiwords make the Urdu word boundary detection a complex and challenging task. This may be a possible explanation for the fact that no standard efficient Urdu word tokenizer is publicly available. An efficient Urdu word tokenization system would be needed to deal with these issues and to properly tokenize Urdu text.

To conclude, the space insertion problem, space omission problem, and ambiguity in tokenizing multi-words make the Urdu word boundary detection a complex and challenging task. This may be a possible explanation for the fact that no standard efficient Urdu word tokenizer is publicly available. An efficient Urdu word tokenization system would be needed to deal with these issues and to properly tokenize Urdu text.

3.2 Challenges for sentence boundary detection

Sentence boundary detection is a nontrivial task for Urdu text because (i) it does not use any special distinguishing characters between upper and lower case, (ii) punctuation markers are not always used as sentence separators, (iii) sentences are written without any punctuation markers, and (iv) there is a lack of standard evaluation and supporting resources. For English and other languages, the difference in upper and lower case is helpful in identifying sentence boundaries. Furthermore, in English language, there is a convention that if a period is followed by a word starting with a capital letter, then it is more likely to be a sentence marker, whereas in Urdu, there are no upper and lower-case distinctions. Punctuation such as “_”, “.”, “?” and “!” are used as sentence terminators, and these can also be used inside the sentence.

The Table 2 shows example sentence boundary markers (SBM) (such as sentences at index i, ii, iii, and iv, in all these sentences question, period, exclamation, and double quotes marker are used at the end of sentences to represent a sentence boundary) and nonsentence boundary markers (NSBM) for Urdu text. It can be observed from these examples that the NSBM are also frequent because they are being used between dates (such as sentence at index vii, in this sentence a period mark is used with in a sentence which is actually not a sentence boundary), abbreviations (index v, this sentence is composed of several period markers; however, first two are not indicating a sentence boundary marker), emphatic declaration (index vi, here exclamation marker is used with in a sentence, that is, not a sentence boundary mark), names and range (index viii, i.e., a first period and double quote marker is used within a sentence, but both are not a sentence ending marker), and sentences without any SBM (index ix). Consequently, these kinds of examples make the sentence tokenization of Urdu text a challenging task.

3.3 Challenges for POS tagging

POS tagging for the Urdu language is also challenging and difficult task due to four main problems (Mukund *et al.* 2010; Naz *et al.* 2012): (i) free word order (general word order is SOV), (ii) polysemous words, (iii) Urdu is highly inflected and morphologically rich, and (iv) the unavailability of gold-standard training/testing data set(s). We briefly discuss these issues here.

First, Urdu sentences have a relatively complex syntactic structure compared to English. Anwar *et al.* (2007b) have shown examples of the free word order and its semantic meaningfulness in the Urdu language. Second, as with other languages, Urdu also has many polysemous words, where a word changes its meaning according to its context. For example, the word باسی (BASY) means “stale” if it is an adjective and “resident” when it is a noun. Third, Urdu is also a highly inflected and a morphologically rich language because gender, case, number, and forms of verbs are expressed by the morphology (Hardie 2003; Sajjad and Schmid 2009). Moreover, Urdu language represents case with a separate character after the head noun of the noun phrase (Sajjad

Table 2. Examples showing Sentence Boundary Markers (SBM) and Non-Sentence Boundary Markers (NSBM) for Urdu text

index	Marker ^a	Text
i	QM-SBM	مشرف کو باہر کون جانی دنا گما ؟ !GYA DYA JANE KYON BAHR KO MSHRF Why was Musharraf let to go abroad?
ii	PM-SBM	انڈیا میں آئی سی سی ورلڈ ٹی ۲۰ کا آغاز - _AAGHAZ KA 20 TY ORLD SY SY AAYY MY ANDYA Inauguration ceremony of ICC world T 20 held in India.
iii	EM-SBM	اس پر بھی عوام نہ سمجھتی تو ! ! TO SMJHY NH AOAM BHY PR IS Even then if public do not understand then!
iv	DQ-SBM	» مہرے خیال میں ان کی وزیر خارجہ ۲۱ اگست کو آرہی ہیں « "HYN RHY AA KO AGST 21 KHARJH OZYR KE AN MY KHYAL MYRE" "In my opinion the foreign minister is visiting on August 21st"
v	PM-NSBM	یو۔ ای۔ اسی میں کافی پاکستانی بستی ہیں - HYN BSTE PAKSTANY KAFY MY AY- AE -YO Many Pakistanis are living in U.A.E.
vi	EM-NSBM	حضور والا ! آپ پوری ملک کی بادشاہ ہیں - - HYN BADSHH KE MLK PORE AAP ! OALA HZOR My lord! You are the king of this country.
vii	PM-NSBM	آج ۲۰۱۶ - ۶ - ۳ بی - - HYN 3-6-2016 AAJ Today is 3rd of May 2015.
viii	PM-NSBM DQ-NSBM	» «پاکستان» ۲ - ۴ سی جیت رہا ہے - - HE RHA JET SE 3 - 2 "PAKSTAN" "Pakistan" is winning by 2-4.
ix	SBM Missing	ابھی چند سال کا وقت لگی گا - ABHEY CHAND SAL KA WAQET LAGEY GA Still this will take few years

^aQM: Question Mark, PM: Period Mark, EM: Exclamation Mark, DQ: Double Quotes

and Schmid 2009). They are sometimes considered as postpositions in Urdu due to their place of occurrence and separate occurrence. If we consider them as case markers, then Urdu has accusative, dative, instrumental, genitive, locative, nominative, and ergative cases (Butt, 1995: p. 10). Usually, a verb phrase contains a main verb, a light verb (which is used to describe the aspect), and a tense verb (describes the tense of the phrase) (Hardie 2003; Sajjad and Schmid 2009). Finally, there is a lack of benchmark training/testing data sets that can be used for the development and evaluation of Urdu POS taggers.

4. Urdu natural language toolkit

This study aims to develop a natural language toolkit for the Urdu language. The UNLT consists of word tokenization, sentence tokenization, and POS tagging modules. The following sections discuss these modules in more detail.

4.1 Urdu word tokenizer

4.1.1 Generating supporting resources for Urdu word tokenizer

For our proposed Urdu word tokenizer, we have developed two dictionaries: (i) a complex words dictionary—to address space insertion problem and (ii) a morpheme dictionary—to address the problem of space omission.

Complex words dictionary: To address the space insertion problem, a large complex words dictionary was created using the UMC Urdu data set (Jawaid, Kamran, and Bojar 2014), which contains data from various domains including Sports, Politics, Blogs, Education, Literature, Entertainment, Science, Technology, Commerce, Health, Law, Business, Showbiz, Fiction, and Weather. From each domain, at least 1000 sentences were randomly selected and preprocessed to remove noise (see Section 5.4). After noise removal, to speed up the dictionary creation process, a basic space-based tokenization approach was implemented in Java to split sentences into words. Space-based tokenization resulted in some incorrect word generation, for example, complex words such as the prefix ان گنت (AN GNT “countless”) is incorrectly split into a morpheme, ان (AN, literally “this”) and a stem, گنت (GNT, literally “count”), postfix حملہ آور (HMLH AAOR, “assailant”) is incorrectly split as حملہ (HMLH, “attack”), that is, a root and آور (AAOR, literally “hour”), that is, a morpheme

Compound words that can be categorized into three types with respect to their formation: (i) *AB formation*—two roots and stems join together, (ii) *A-o-B formation*—two stems or roots are linked to each other with the help of و (wao) (a linking morpheme), and (iii) *A-e-B formation*—“e” is the linking morpheme which shows relation between A and B (for more detailed discussion see Rehman, Anwar, and Bajwa 2011). In this research, all three types have been used without any classification, for example, A-o-B formation type of compound word غور و فکر (GHOR O FKR, “contemplation”) is incorrectly split as غور (GHOR, literally “ponder”) a root, و (O) a linking morpheme, and a stem فکر (FKR, literally “worry”). Reduplication that have two types: (i) *full reduplicated word*—two duplicate words are used to form a word and (ii) *echo reduplication*—the onset of the content word is replaced with another consonant (detailed information can be found in Bögel et al. 2007). Echo reduplication word, دن بدن (DN BDN, “day by day”) is incorrectly split as دن (DN, literally “day”), that is, content word and بدن (BDN, literally “body”), a consonant. One million space-based tokenized words list (henceforth UMC-Words) has been used to form a large complex words dictionary containing: (i) affixes, (ii) reduplications, (iii) proper nouns, (iv) English words, and (v) compound words.

To collect affixes (prefixes and postfixes) complex words from the UMC-Words list (Jawaid et al. 2014), a two-step approach is used. In the first step, a list of prefixes and postfixes are manually generated. In the second step, an automatic routine is used to extract words containing affixes from the large UMC-Words list. Using prefixes and postfixes, the previous and next words are extracted, respectively, from the UMC-Words list.

Reduplications complex words are collected using two methods: (i) full extraction and (ii) partial extraction. The full extraction method is used to extract the full reduplicated words such as جسی جسی (JYSY JYSY, “as”). To extract such full reduplicated words, we compared each word in the UMC-Words list to the next word, if both are the same then concatenate both to form a full reduplicated compound word. The partial extraction method is used to collect the words of echo reduplication, that is, in which a consonant word is a single edit distance away from the first

content word. The echo reduplication words can be further collected using two methods: (i) one insertion extraction and (ii) single substitution extraction.

One insertion extraction method extracts the one insertion reduplicated words, in which the consonant word has one insertion in its content word, for example, دن بدن (DN BDN, “day by day”). It can be noted that the consonant word بدن (BDN, literally meaning “body”) has one more character (three) as compared to the content word دن (DN, literally “day”) (which have two characters). Furthermore, the last two characters of the consonant word are identical to the content word. To extract one insertion reduplicated words, we used the UMC-Words list. The extraction process works as follows: after excluding the first character, if the remaining characters of consonant word are identical as well as having the same character count to the content word, they are one insertion reduplicated word(s) we concatenated them to form a single word.

The single substitution extraction method extracts the single substituted reduplicated word(s)—here the consonant word has single substitution in its content word, for example, خلط ملط (KHLT MLT, “intermixed”). It is worth noting that both words content خلط (KHLT, literally “bad”) and consonant ملط (transliteration: MLT) has three characters and the final two characters are overlapping. To extract one substituted reduplicated word(s) we used automatic routine and applied the following process over the UMC-Words list as: if the length of the content word is matched with the length of the consonant word and the length of content word is greater than two^P characters, and if one character is dissimilar after comparing character by character, then it will form a single substitution reduplicated complex word.

To automatically extract abbreviations (91) and proper nouns (2K), regular expressions are used and further supplemented by manual checking to increase the size of the proper nouns (3K) and abbreviations lists (187). The remaining 65K proper noun list was generated in another NLP project and are used in this study for Urdu word tokenization. In addition to this, manual work^Q was also carried out to remove noisy affix entries. Moreover, compound words (of formation AB and A-e-B) and English words are added to increase the size of the complex words dictionary. However, to collect words of A-o-B formation automatically, a linking morpheme (و, O) has been used. While using a linking morpheme both previous and next words are extracted from the UMC-Words list to form a A-o-B compound words.

The complete statistics of the complex words dictionary are as follows: there are in total 80,278 complex words (7820 are affixes, 278 are abbreviations, 10,000 are MWEs, 1480 are English words, 60,000 are proper nouns, and 700 are reduplication words).

Morpheme segmentation process: To address the space omission problem (see Section 4.1.2), a large-scale morpheme dictionary is automatically compiled from the HC data set (Christensen 2014). Before we proceed further toward the approach used to generate the morphemes dictionary, it is worth describing the morpheme types. Urdu language morphemes can be categorized into (i) free and (ii) bound morphemes. Proposed word tokenizer has to rely on both categories. The bound or functional morphemes such as affixes include prefixes, for example, “گ، لا، کو” (GA, LA, KO), linking morphemes, for example, و، ا (A, O) or suffixes, for example, دہ، شدہ (transliteration: SHDH, ZDH), can only expose their meanings if they are attached to other words, that is, they cannot stand alone. Whereas free or lexical morphemes can stand alone, for example, غم، چست، علم، مقبول (MKBOL, CHST, ALM, GHM, “grief, knowledge, clever, famous”).

^PTo make sure the two character words or auxiliaries could not be erroneously identified as reduplication such as کر کے (KR KE, literally “by doing”).

^QFive undergraduate NLP students have been employed to carry out manual tasks, all are native Urdu speakers and have an interest in Urdu NLP and literature. Furthermore, each student undertook a practical training session on annotation tasks. Each student was given an annotation assignment of 80 random sentences from the UPC data set and requested to extract affixes, compound words, abbreviations, and English words. These assignments were marked and each student was awarded with a score. Students having scored 85% or above were thus selected for annotation tasks.

There are two further categories of free morphemes: (i) true free morphemes and (ii) pseudo-free morphemes. True free morphemes can be either standalone (for e.g., دل (DL, “heart”)) or form part of other words (e.g., درد دل (DRD DL, “angina pectoris”). Pseudo-free morphemes can be a character, affix, or word.

The preceding discussion summarizes the various types of morphemes. However, from a computational linguistics view, free and bound morphemes play a vital role in Urdu word formations (Khan *et al.* 2012); hence, they will be used without any further classification in our proposed UNLT word tokenizer module.

In order to generate the morpheme dictionary, the 1000 most frequent words which have more than 20 occurrences in the HC data set are used (Christensen 2014); the selected words were split to form a morpheme dictionary. The whole chopping process is completed in two steps: (i) Crude-Morphemes (CM) chopping and (ii) Ultra-Crude-Morphemes (UCM) chopping.

In the first step, the first n character(s) of each word are kept while the rest are discarded. For example, in case of $n = 1$, we kept only the first character and discarded all others, thus words such as واقفت (OAKFYT, “awareness”) will return و (wao). Such single character morphemes are helpful to formulate compound words, for instance خوش و خرم (KHSO O KHRM, “canty”). Furthermore, we keep chopping all the words repeatedly with the following values of $n = 2, 3, 4, 5, 6$.^f This process returns واقفت، واقفی، واقف، واق، وا (transliterations are: OA, OAK, OAKF, OAKFY, OAKFYT) morphemes for the word واقفت (OAKFYT, “awareness”). There may be a situation where we may lose several valuable morpheme(s), if the length of $n > 6$. Nevertheless, this is a rare case. Henceforth, we will call this method Crude-Morpheme chopping.

To generate entirely different morphemes from the same word, we further applied a modified version of CM chopping, that is, UCM. In which, we skipped the first character and then applied the CM chopping with length $n = 2, 3, 4$. Thus, UCM chopping resulted with these morphemes، اقفت، اقفی، اق، اق (transliterations are: AK, AKF, AKFY, AKFYT) for the word واقفت (OAKFYT, “awareness”). Furthermore, we iterate the UCM chopping method by skipping the first two characters (as well as three, four etc.), until we meet the last two characters. Thus, the following morphemes are returned by UCM, in the third قفت، قفی، قف (transliterations are: KF, KFY, KFYT), in the fourth فت، فی (transliterations are: FY, FYT), and in the last بت (transliteration, YT) iterations.

Repeating CM and UCM chopping on the entire list of words will return all possible morphemes. The two chopping methods used in this study will result in erroneous morphemes. However, we manually examined the morpheme dictionary and removed these. The number of morphemes generated by the CM and UCM chopping methods is 5089 and 7376, respectively.

It can be observed from the above discussion that two different large-scale dictionaries, that is, the complex words dictionary and the morphemes dictionary are generated with distinct approaches and with various statistics. These dictionaries will be used to solve the space omission and space insertion problems with the word tokenizer module of UNLT. To the best of our knowledge, no such large complex words (a study Hautli and Sulger 2009 just proposed a scheme to extract location and person name), and morpheme dictionaries have been previously compiled semi-automatically for Urdu, to perform Urdu word tokenization.

4.1.2 Proposed Urdu word tokenizer

To investigate an effective approach for UNLT Word Tokenization (henceforth UNLT-WT approach), our method (see Algorithm 1) is a combination of state-of-the-art approaches: rule-based maximum matching, dictionary lookup, statistical *tri*-gram Maximum Likelihood Estimation (MLE) with back-off to *bi*-gram MLE. Furthermore, smoothing is applied to avoid data

^fAn assumption made by us after analysing Urdu text that a word is formed of a maximum of six morphemes.

Algorithm 1 UNLT-WT approach

Step 1: Initialize $flag_bit=false$, $row=1$, $column=1$, $word_counter=0$;

Step 2: Create array $words_list[row][column]$, array $morphemes_list$, array $compound_words_list$;

Step 3: Remove all white spaces and ZWNJ, to form a space free input text.

Step 4: Read *bi*-gram of input text.

Step 5: Match this *bi*-gram with each word of *morphemes_list*

Step 6: Extract all those morphemes from *morphemes_list*, which matched with *bi*-gram.

Step 7: Store each extracted morpheme on a separate row/column of *words_list*

Step 7.1: For each row, copy the $flag_bit$, $word_counter++$

Step 8: If no match is found in *morphemes_list*, split *bi*-gram into *uni*-gram.

Step 8.1: Store the first *uni*-gram with previous morpheme (column) except ๑ (character O) and ๑ (character A) (use in compound words) and turn the $flag_bit=true$. For ๑ and ๑, store it on separate column of array $words_list[row][column]$ and increment $word_counter$.

Step 9: Repeat the steps 4 to 8, until sentence ending marker, and for each row of *words_list*.

Step 10: Select the row having minimum $word_counter$ value and $flag_bit=false$.

Step 11: If multiple rows are qualified in step 10 then

Step 11.1: Calculate *tri*-gram MLE for each row.

Step 11.1.1: Select the one having highest value of *tri*-gram MLE.

Step 11.2: If in step 11.1.1, any row having *tri*-gram MLE value equal to *Zero*, then calculate *bi*-gram MLE for each row.

Step 11.2.1: Select the one having highest value of *bi*-gram MLE.

Step 11.3: If in step 11.2.1, any row having *bi*-gram MLE value equal to *Zero* then, calculate *bi*-gram *smoothing* for each row.

Step 11.3.1: Select the one having highest value of *smoothing*.

Step 12: For final selected row, read each column and match in the compound word dictionary.

Step 12.1: If a match is found, then read the next column of selected row in step 12 and repeat step 12 for the remaining part of selected row.

Step 12.1.1: If complete match is found then concatenate with the columns in step 12.1.

Step 12.1.2: Move each element of final selected row in step 12, decrease the array index.

Step 13: Finally, list of tokenized word will be produced.

sparseness. A step-by-step working example of the proposed algorithm can be seen at.⁸ However, this section just presents the statistical approach used to solve space omission problem.

Maximum likelihood and smoothing estimation: In our proposed UNLT-WT approach (see Algorithm 1) at step 11.1, we used *tri*-gram MLE and smoothing estimations because there can be multiple tokenized sequences for which $flag_bit=false$ and $word_count$ are equal. For instance,

⁸<https://doi.org/10.17635/lancaster/thesis/831> - pp. 88–98.

there are two given texts, (i) اسي باهر جا كي پڙهني دو (transliteration: ASE BAHR JA KE PRHNE DO, “let him go abroad for higher studies”) and (ii) اسي باهر جي كي پڙهني دو (transliteration: ASE BAHR JY KE PRHNE DO, literally meaning “let him yes abroad for higher studies”). Both have six tokens with flag_bit= false, but only the first text is semantically correct and meaningful. For such ambiguous cases, we calculate an N -gram language model with MLE for parameter and Laplace for smoothing estimation. The goal of these estimations is to find an optimized segmented sequence with the highest probability. This can be shown by a given mathematical expression, a general statistical model of our proposed UNLT-WT approach.

$$\hat{P}(m_j|m_{1-N+1}m_{n-N+2}...m_{j-1}) = \arg \max_{M \in \zeta(I|D)} P(M) \tag{1}$$

Here, $\zeta(I|D)$ denotes all possible tokenized words of the input string, that is, $I = i_1i_2...i_l$ with l characters, and M denotes string concatenation of all possible tokenized sequences, that is, $M = m_1m_2...m_n$, in terms of morphemes dictionary D . Theoretically, it is assumed that the n -gram model outperforms with a high value of N . However, practically the data sparseness restricts better performance with high order N . Therefore, in our UNLT-WT approach, we opted for *tri*-gram ($N = 3$) or *bi*-gram ($N = 2$) MLE. These have proved to be successful in several tasks for resolving ambiguity (e.g., POS tagging Brants 2000, automatic speech recognition Abdelhamid, Abdulla, and MacDonald 2012 and word tokenization Fu, Kit, and Webster 2008).

The task of resolving similar sequence ambiguities for the above two texts is accomplished by using *tri*-gram MLE (Jurafsky and Martin 2014) as

$$P(t_j|t_{j-2}, t_{j-1}) = \frac{C(t_{j-2}, t_{j-1}, t)}{C(t_{j-2}, t_{j-1})} \tag{2}$$

where t represents the individual token, C is a count of three ($t_{j-2}t_{j-1}t$) and two ($t_{j-2}t_{j-1}$) consecutive words in the data set, and P is the *tri*-gram contestant MLE value of each of the possible segmented sequences. The calculated probability for the first sequence is 3.2e-08 while for the second it is 0.

As *tri*-grams take account of more context, if this specific context is not found in the training data (see Section 5.1), we back-off to a narrower contextual *bi*-gram language model. *Bi*-gram cumulative probability values have been calculated as given by Jurafsky and Martin (2014):

$$P(t_j|t_{j-1}) = \frac{C(t_{j-1}t)}{C(t_{j-1})} \tag{3}$$

where t represents the individual token, C is a count of two ($t_{j-1}t$) and one (t_{j-1}) consecutive word(s) in the data set, and P is the *bi*-gram contestant MLE value of each of the possible segmented sequences. The calculated probability for the first sequence is 2.7e-6 for the former sequence and 0 for the latter one.

These *zero* probabilities are again an underestimation of the input string, ultimately a cause for the data sparseness. Even if a statistical language model is trained on a very large data set, it will remain sparse in some cases. However, there is always a possibility that the input text occurs in the test data set (Chen and Goodman 1999), thus assigning them to zero made this an unstable, frail, and specific estimator. Therefore, to overcome this, different *smoothing* techniques have been proposed in previous literature (Jurafsky and Martin 2014) with different characteristics (such as smoothing the probability, etc.). Hence, it is primarily aimed at making a robust and generalize language model by re-evaluating lower or zero probability upward and vice-versa for high probabilities.

For this study, we employed Laplace (add-one) smoothing (Jeffreys 1998), as one of the oldest, simplest, and baseline estimations. This estimation adds one to all frequency counts, that is, that

all *bi*-gram probability counts have been seen one more time than actually exists in the training data as

$$P_{add:1}(t_j|t_{j-1}) = \frac{1 + C(t_{j-1}, t)}{V + C(t_{j-1})} \quad (4)$$

where v represents the unique words (types), added to the total number of words $C(t_{j-1})$ in order to keep the probability normalized (Jurafsky and Martin 2014). We have used Laplace smoothing to estimate the parameters required for data sparseness in order to increase the *bi*-gram MLE value for *اسی باہر جی کی پڑھنی دو* (transliteration: ASE BAHR JY KE PRHNE DO, “let him go abroad for higher studies”), from 0 to $1.9e-14$, and decreased value for *اسی باہر جا کی پڑھنی دو* (transliteration: ASE BAHR JA KE PRHNE DO, literally meaning “let him yes abroad for higher studies”), from $2.7e-6$ to $3.8e-7$. As the latter tokenized sequence has the highest smoothing MLE. Therefore, it will be selected by UNLT-WT as the best tokenized sequence, which is correct.

4.2 Urdu sentence tokenizer

4.2.1 Rule-based approach

For our proposed rule-based approach, to manually extract rules for the sentence tokenization task, initially, a subset of the UMC data set (Jawaid *et al.* 2014) composed of 13K sentences is selected, which contains Urdu text from various domains or genres including News, Religion, Blogs, Literature, Science, and Education. After preprocessing (see Section 5.4), we retained 10K sentences, which were used to extract rules to develop our proposed Urdu sentence tokenizer. The rules were devised to include sentence termination markers (., ?, » and !), regular expressions, and supplementary dictionary lookup[†] (henceforth UNLT-ST-RB approach). These heuristics are applied as follows:

1. If the current character is a period marker (.) AND the same mark appears after two or three characters, then consider it as an abbreviation and match it in the abbreviation list.
2. If within the next 9 characters (from any previous SBM marker), an exclamation mark (!) is found, then this is not a sentence boundary marker.
3. If the character before a double quote (») is a period (.) or question (?) mark, then it is a sentence boundary marker.
4. Apply regular expressions for detecting the date and hyphenated numeric values.
5. In addition to this all the above rules from 1 to 4, split sentences based on the question (?), period (.), and exclamation (!) markers.

4.2.2 ML-based approach

In this approach, we are exploiting a support vector machines (SVM) classifier (Hearst *et al.* 1998) to detect the sentence boundaries of the Urdu text—using the features described below another approach is formed, that is, UNLT-ST-ML. SVMs offer robust classification even with sparse vectors of large dimension (Akita *et al.* 2006), its good performance results on textual data and its suitability for binary classification (Kreuzthaler *et al.* 2015) task make this a suitable classifier for sentence boundary detection. Moreover, SVMs use a function (see Equation (5)) for classifying sentence boundary label pairs (x_j, y_j) , $j = 1, \dots, m$ for all $x_j \in \mathbb{R}^n$ to a target value $y \in \{1, -1\}$. Where

[†]We used the same dictionary compiled for the word tokenization task (see Section 4.1.1)

$w \in \mathbb{R}^n$ a weight coefficient and $b \in \mathbb{R}$ is a bias. We are using a polynomial kernel implemented in Weka.^u

$$f(X) = \text{sgn}(w^T \phi(x) + b) \quad (5)$$

Features for ML approach:

- Probability (UMC data set^v Jawaid *et al.* 2014 is used) that a word with “؟, _ and !” occurs at the end of a sentence
- Probability (UMC data set^w is used Jawaid *et al.* 2014) that a word with “؟, _ and !” occurs at the beginning of a sentence
- Length of a word with “؟, _ and !”
- Length of a word after “؟, _ and !”
- Is a sentence contains an abbreviation
- Is a sentence contains a date/numeration
- *Bi-* and *tri-*grams words information (preceding “؟, _ and !”) are used
- If a word before “؟, _ and !” markers contains any one of the tag (NN, NNP, JJ, SC, PDM, PRS, CD, OD, FR, Q, and CC. See Section 4.3.1 for POS tags) is not a sentence boundary

4.3 Urdu part of speech tagging

4.3.1 Existing Urdu POS tagset

The tagging accuracy of a POS tagger is not only dependent on the quality and amount of training data set but also on the POS tagset used for annotation. In the prior literature, we found three commonly used POS tagsets for the Urdu language:

- (i) Hardie’s POS tagset (Hardie 2004),
- (ii) Sajjad’s POS tagset (Sajjad 2007), and
- (iii) Centre for Language Engineering (CLE) Urdu POS tagset (Ahmed *et al.* 2014).

Hardie’s POS tagset (Hardie 2004) was an early attempt to resolve the grammatical tag disambiguation problem for the Urdu language. This tagset follows the EAGLES^x guidelines and consists of 350 morphosyntactic tags, which are divided into 13 main categories. Some grammarians (Platts 1909) propose only three main categories, whereas (Schmidt 1999) used 10 main categories for Urdu text. There were a number of shortcomings observed in Hardie’s POS tagset (Hardie 2004). For example, the possessive pronouns like *میرا* (MYRA “my”), *تمہارا* (TMHARA “your”), and *ہمارا* (HMARA “our”) are assigned to the category of *possessive adjective*, which is incorrect. Many grammarians marked them as *pronouns* (Platts 1909; Javed 1985). Moreover, the Urdu language has no *articles*, but this tagset defined articles. Another issue with the tagset is the use of *locative* and *temporal adverbs* such as *ہاں* (YHAN “here”), *وہاں* (OHAN “there”), and *اب* (AB “now”), which are treated as *pronouns*. The *locative* and *temporal* nouns such as *صبح* (SBH “morning”), *شام* (SHAM “evening”), and *گھر* (GHR “home”) appear in a very similar syntactic context. To conclude, these grammatical misclassifications as well as the large number of POS tags with relatively small training data will affect the accuracy of POS taggers developed for the Urdu language.

^u<http://www.cs.waikato.ac.nz/ml/weka/> - last checked: 22-September-2020.

^vTo calculate that a certain word occurs before a sentence boundary.

^wTo calculate that a certain word occurs after a sentence boundary.

^x<http://www.ilc.cnr.it/EAGLES96/home.html> - Last checked: 07-December-2016.

Another POS tagset (henceforth Sajjad’s POS tagset) (Sajjad and Schmid 2009), consists of 42 POS tags with finer grained categories for pronouns and demonstratives. However, it is lacking in terms of Urdu verb, tense, and aspect.

A recently released CLE Urdu POS tagset (Ahmed *et al.* 2014) contains 35 tags and addresses most of the issues reported above. It is based on the critical analysis of several previous iterations of Urdu POS tagsets. Furthermore, it is built on the guidelines of the Penn Treebank^y and a POS tagset for common Indian languages.^z In the CLE Urdu POS tagset, a verb category has multiple tags based on the morphology of the verbs. Furthermore, it has shown promising results on Urdu text see Section 2.3).

For this study, we selected the CLE Urdu POS tagset (Ahmed *et al.* 2014) for the following reasons: (i) it provides correct grammatical classifications, (ii) it provides purely syntactic categories for major word classes, and (iii) provides reasonable performance on a small size test data set.

4.3.2 Proposed Urdu POS tagging approaches

For this study, we applied two stochastic approaches for Urdu POS tagging: (i) *tri*-gram Hidden Markov Model and (ii) Maximum Entropy-based model. The reason for selecting these two methods for Urdu POS tagging is many fold, (a) they have proven to be effective for POS tagging not just for English (Yi 2015) but also for other languages which are closely related to Urdu such as Hindi (Dandapat 2008; Joshi, Darbari, and Mathur 2013), (b) both are well established stochastic models for automatic POS tagging task (Wicaksono and Purwarianti 2010), (c) these methods have been primarily investigated for when dealing with languages with limited resources (Azimzadeh, Arab, and Quchani 2008; Ekbal, Haque, and Bandyopadhyay 2008), and (d) these models have not been previously compared for the Urdu language.

Hidden Markov Model (HMM) for POS tagging: In general, the Urdu POS tagging task can be formulated as: given a sequence of words w_1, \dots, w_n , find the sequence of POS tags t_1, \dots, t_n from a POS tagset^{aa} using some statistical model. In this section, we have used HMM stochastic learning model described by Rabiner (1989), while Thede and Harper (1999) redefined it for the POS disambiguation task. This model was implemented in Garside and Smith (1997), Bird *et al.* (2009) for POS tagging. For our experiments, we used a third-order HMM learning model, also referred to as a *tri*-gram POS tagging. This model is composed of transitional (contextual) and lexical (emission) probabilities as:

$$\hat{T} = \arg \max_{t \in T} P(t_1, \dots, t_n | (w_1, \dots, w_n)) \tag{6}$$

Using Bayes’ theorem, the above equation can be rewritten as for 3rd order model as

$$P(t_1, \dots, t_n | w_1, \dots, w_n) = \arg \max_{t_1, \dots, t_n} \prod_{j=3}^n \underbrace{(P(t_j | t_{j-1}, t_{j-2}))}_{TransitionProbability} * \prod_{i=1}^n \underbrace{P(w_i | t_i)}_{LexicalProbability} \tag{7}$$

During the training process, the above *tri*-gram HMM language model computes two probability factors for the sequences: (i) emission probabilities, aimed at determining the probability of a particular tag conditioned on particular word and (ii) transitional probabilities, used to find the probability of a particular tag on the basis of given preceding tag(s). Given a sentence, the aim of the HMM language model is to search the tagging sequence and choose the most likely sequence

^y<https://www.cis.upenn.edu/treebank/> - Last checked: 13-June-2017.

^z<https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/I08-7013.pdf> - Last checked: 24-September-2016.

^{aa}35 tags as in CLE Urdu POS tagset: <http://www.cle.org.pk/software/langproc/POSTagset.htm> - Last checked: 20-November-2016.

that maximizes the dot product of lexical and transition probabilities. That can be computed by using a Viterbi algorithm (Viterbi 1967).

Parameter estimation: We can estimate the HMM parameters by applying the simplest *tri*-gram MLE (see Section 4.1.2), used for computing relative frequencies. We have used a training data set (see Section 5.3) to find tag frequency counts (C) for two or three consecutive tag pairs (t_{j-2}, t_{j-1}, t_j) , (t_{j-2}, t_{j-1}) . Where t_j is the j th tag of annotated data set used during training process. The following equation requires frequencies count of $w_i t_i$, where w_i is the word and t_i is the tag assigned to i th word. The *tri*-gram language model (see Section 4.1.2) and the following equation is used with these parameter settings, $1 \leq (i, j) \leq n$.

$$P(w_i | t_i) = \frac{C(w_i t_i)}{C(t_i)} \quad (8)$$

Smoothing: We have used the MLE for parameter estimation (see Section 4.3.2), consequently, such models may come across a situation where unseen events do not occur or have quite low frequencies in the trained model. Therefore, the zero probability of such occurrences produces problems in the multiplication of probabilities, eventually, leading to a data sparseness.

To avoid data sparseness, we need some estimators that automatically assign a part of the probability mass to unknown words and tag sequences, thus yielding an improvement for unseen events and overall accuracy improvement for the POS tagger. For this, different smoothing techniques have been cited in the literature with an objective to decrease the probability of seen events and assigning appropriate nonzero probability mass to unseen events. In this study, five different smoothing techniques were adopted including (i) linear interpolation, (ii) Laplace, (iii) Lidstone's, (iv) Good-Turing, and (v) Kneser-Ney estimations. Adopting them with an HMM model thus alleviates sparse data issues.

Linear interpolation: A well-practised smoothing technique consists of linearly combined estimation for different order n -grams as

$$P(t_i | t_{i-1}, t_{i-2}) = \lambda_1 \rho(t_i) + \lambda_2 \rho(t_i | t_{i-1}) + \lambda_3 \rho(t_i | t_{i-1}, t_{i-2}) \quad (9)$$

where P is a valid probability distribution, ρ are maximum likelihood estimates of the probabilities, and $\lambda_1 + \lambda_2 + \lambda_3 = 1$ to normalize the probability. Although there are different ways to estimate λ s, for our experiments, we adopted *deleted linear interpolation*, cited in Brants (2000).

The deleted linear interpolation successively removes each *tri*-gram from the training data set. Moreover, this technique estimates the best value for the λ s from all other n -grams in the data set, making sure that the value of λ does not depend upon the particular n -gram. Further, it computes the weights depending on the counts of each i -gram, involved in the interpolation. Thus, our first HMM-based proposed model is a combination of linear interpolation smoothing technique along with *tri*-gram HMM model (henceforth T-HMM-LI).

Laplace and Lidstone's estimation: Laplace estimation (one of the oldest and simplest smoothing techniques) updates the count by one of each bi -gram occurs compared to the actual frequency in training data (Jurafsky and Martin 2014) (see Section 4.1.2). Whereas Lidstone's smoothing estimation (Manning and Schütze 1999) generalizes Laplace, by adding an arbitrary value to all (seen or unseen) the events. Although the values for λ can be calculated using different methods, for our experiments we used the same value cited in the research article (Manning and Schütze 1999), that is, a well-known expected likelihood estimation (ELE). Thus, Lidstone's estimation (Manning and Schütze 1999) can be calculated as

$$P_{\text{lidstone}(x,\lambda)} = \frac{\lambda + C(X)}{V\lambda + N} \quad \lambda = 0.5 \quad (10)$$

where V represents the unique words (vocabulary) against the total number of words N to keep probabilities normalized (Jurafsky and Martin 2014). The generalized formulation of Lidstones and Laplace estimation in an HMM-based Urdu tagger is as follows:

$$\pi_i = \frac{C(s_i(t=0)) + \lambda}{C(tokens) + V_{tag}\lambda} \tag{11}$$

$$a_{ij} = \frac{C(s_i \rightarrow s_j) + \lambda}{C(tokens) + V_{tag}\lambda} \tag{12}$$

$$P(s_j) = \frac{C(s_j) + \lambda}{C(tokens) + V_{tag}\lambda} \tag{13}$$

$$P(w_k) = \frac{C(w_k) + \lambda}{C(tokens) + V_w\lambda} \tag{14}$$

Here, V_{tag} is the number of possible tags and V_w is the size of the approximated vocabulary. Our second POS tagging model is a combination of Laplace and *tri*-gram HMM model (henceforth T-HMM-LaE). The third POS tagger makes use of Lidstone’s estimation and supervised *tri*-gram HMM model parameters (we shall call this T-HMM-LiE).

Good–Turing estimation: Uses the probability mass of n -grams (that occur $c + 1$ times) which is seen once to re-estimate the count of n -grams (that are seen exactly c times) that were never seen. This can be described as

$$c^* = (c + 1) \frac{n_{c+1}}{n_c} \tag{15}$$

The fourth POS tagger makes use of Good–Turing estimation and supervised *tri*-gram HMM model parameters (hereafter, T-HMM-GT).

Kneser–Ney smoothing: This outperforms all other smoothing techniques. In this paper, we have used the modified version of this smoothing (Christer 1996), which is an interpolated variation of Kneser–Ney smoothing with an augmented version of absolute discounting, thus the transition probabilities $p(t_j|t_{j-1})$ are calculated as

$$P(t_j|t_{j-1}) = \frac{f(t_{j-1}t_j) - D(f(t_{j-1}t_j))}{\sum_{t_j} f(t_{j-1}t_j)} + \gamma(t_{j-1}) \cdot \frac{N_{1+(\cdot, t_i)}}{N_{1+(\cdot \cdot)}} \tag{16}$$

Here D is estimated value and $D(f) = \{0 \text{ if } f = 0, D_1 \text{ if } f = 1, D_2 \text{ if } f = 2, \dots\}$. Another POS tagger is formed using supervised *tri*-gram HMM model and Kneser–Ney smoothing (henceforth T-HMM-KN).

Maximum Entropy (MaEn) model for POS tagging: The other adopted stochastic learning model is MaEn, and we aimed to compare this to the above described *tri*-gram HMM-based models, to find the most optimal POS tagger for Urdu. The MaEn statistical assumption is a simplistic model, it assigns a probability distribution for every tag, given a word and its context as

$$P(t_1, \dots, t_n | w_1, \dots, w_n) = \prod_{j=1}^n P(t_j | c_j) \tag{17}$$

where t is the individual tag in the set T of all possible tags, that is, t_1, \dots, t_n for a given a sentence, c is defined as the context, usually defined as the sequence of words w_1, \dots, w_n and the tag preceding

the word. The maximum likelihood tag sequence is used for assigning probabilities to a string of input words.

The principle of estimating probabilities in MaEn model is to make as few assumptions as possible, other than the constraint imposed. Furthermore, these constraints are learned from the training data, which express some relation between features extracted and outcome. The probability distribution which satisfies the above property has the highest entropy; thus, it agrees with the maximum likelihood distribution and has a general form as cited in Ratnaparkhi (1996):

$$P(t|c) = \frac{1}{N} \exp \sum_{j=1}^k \alpha_j f_j(c, t) \quad (18)$$

where N is the total number of training samples (normalization constant), and f_j is feature function on the event (c, t) . Feature functions used by MaEn model are binary valued and defined to capture relevant aspects of language. The α_j is a model parameter with k features, which is determined through the Generalize Iterative Scaling (GIS) algorithm (Curran and Clark 2003). However, these model values and features are primary ingredients of MaEn learning model.

Features selection in MaEn model: As described previously, the MaEn is feature-based probabilistic model, to obtain high accuracy we used two binary valued features that might be helpful for predicting POS tag, these are determined empirically for Urdu POS tagging along with MaEn model as (i) context window and (ii) word number. The best context window with five words has been identified, which is composed of n -gram $(W_{i-2}, W_{i-1}, W_i, W_{i+1}, W_{i+2})$ and n -POS $(t_{i-2}, t_{i-1}, \text{ and } t_i)$ information. If the current word is a number such as “۲ ۵ ۰ ۳ ۱”, another feature can be created:

$$f_j(c, t) = \begin{cases} 1 \text{ if } \text{WordReadIsNumber}(w_j) = \text{true and } t_j = \text{CD} \\ \text{else } 0 \end{cases} \quad (19)$$

Using the above-mentioned features with MaEn, we formulated another Urdu POS tagging model (henceforth MEn). However, these suitable binary valued features are the same for other languages. We examine some other important feature sets for the Urdu language below.

Morphological information for HMM and MaEn models: To improve the tagging accuracy of the unknown words in the above models, we have developed an exclusive feature set after detailed analysis of the UNLT-POS training data set (see Section 5.3). This feature set is intended to have the capability to capture lexical and morphological characteristics (features) of the Urdu language. The captured morphological features are based on information retrieved from a stemmer^{bb} and dictionary,^{cc} assuming that information is complete.^{dd} Thus, we boosted the lexical probability of assigning restricted lexical (POS) tag to a word. Consequently, the integrated models are expected to perform better with such artificial weight (reduced set of possibilities) for a given word. All the above models (T-HMM-LI, T-HMM-LaE, T-HMM-LiE, T-HMM-GT-MA, T-HMM-KN, and MEn) are incorporated with such restricted POS tags features, henceforth, T-HMM-LI-MA, T-HMM-LaE-MA, T-HMM-LIE-MA, T-HMM-GT-MA, T-HMM-KN-MA, and MEn-MA.

The above-mentioned MA information is helpful to restrict the possible choice of POS tags for a given word, on the other hand, suffix and prefix (of current word) information can also help us to further improve the POS models. For HMM-based POS models, suffix information has been used during the smoothing of emission probabilities. Whereas for the MEn model, the suffix and

^{bb}<http://www.cle.org.pk/software/langproc/UrduStemmer.htm> - Last checked: 17-November-2016.

^{cc}<http://182.180.102.251:8081/oud/default.aspx> - Last checked: 09-October-2016.

^{dd}If a word is unknown, then it belongs to one of the open class lexical categories, that is, all classes of Noun, Adjective, Verb, Adverb, and Interjection.

prefix information are used as another type of feature. It is extended using a prefix and suffix up to a length of four. It is also important to note, using prefix and suffixes of length ≤ 4 for all words in MEn gives better results instead of using only rare words as described by Ratnaparkhi (1996). The primary reason for much improved results based on prefix and suffix is that a significant number of instances are not found for most of the word of the language vocabulary, with a small amount of annotated data. HMM based (T-HMM-LI, T-HMM-LaE, T-HMM-LiE, T-HMM-GT, and T-HMM-KN), and MEn models are incorporated with suffix information, we shall call them T-HMM-LI-Suf, T-HMM-LaE-Suf, T-HMM-LiE-Suf, T-HMM-GT-Suf, T-HMM-KN-Suf, and MEn-Suf POS taggers.

The last six POS models represent combinations of various statistical, smoothing, and features as described above. The T-HMM-LI-Suf-MA is a combination of *tri*-gram HMM along with Linear interpolation, restricted POS tags feature, and suffix information. T-HMM-LaE-Suf-MA is based on the *tri*-gram HMM model with further incorporation of Laplace smoothing, suffix, and restricted POS tags. In T-HMM-LiE-Suf-MA, we have used *tri*-gram HMM along with Lidstone's estimation, with suffix, and restricted POS tags. The T-HMM-GT-Suf-MA is a combination of *tri*-gram HMM along with Good-Turing, restricted POS tags, and suffix information. In T-HMM-KN-Suf-MA, *tri*-gram HMM, Kneser–Ney estimation, suffix, and restricted POS tags are used. MEn-Suf-MA POS tagging model is a collection of, MaEn, contextual window, suffix, and restricted POS tags.

4.4 Deep learning approaches for comparison

We have performed an empirical comparison of our proposed word/sentence tokenizers and POS tagger, with the recently proposed Trankit NLP Toolkit (Nguyen *et al.* 2021). We have selected this deep learning toolkit for the evaluation process as this has reported good results as compared to other deep learning methods. This is a light-weight transformers-based (a deep neural network-based new SOTA architecture in NLP Vaswani *et al.* 2017) tool for multi/monolingual text processing, with trainable and pretrained pipelines for NLP tasks for 100 languages including Urdu. Trankit's recent version^{ee} NLP tools are built upon the transformer XLM-RoBERTa (Conneau *et al.* 2020) (a cross-lingual transformers based masked language model, trained on more than 2 terabytes of preprocessed Common Crawl data for 100 languages), extends the state-of-the-art in word/sentence tokenization, dependency parsing, part-of-speech tagging, and tagging of morphological features. It also attained high performance in various other basic NLP tasks including word tokenization, multiword expansion, and lemmatization of the 90 treebanks.

Trankit's word and sentence tokenizers are based on word parts rather than character-based to exploit the contextual information (Kudo 2018). It begins on a given input text t , by splitting it into substrings on the basis of space. Furthermore, a multilingual sentence splitter (Kudo and Richardson 2018) further breaks each substring into word pieces, concatenated to obtain an overall sequence of word parts $w = [w_1, w_2, w_3, \dots, w_n]$ for t . In the next step, a pretrained transformer takes w as input to produce corresponding representation vectors $v = [v_1, v_2, v_3, \dots, v_n]$ for each word parts in the sequence w . Each vector from v will be fed to a feed forward neural network with a *softmax* at the end to predict whether the w_i is the end of a single word, multiword, or a sentence. Finally, all of the predictions for all word parts w are accumulated to decide a single-word token, multiword token, and sentence boundaries for the given input text t . Word and sentence tokenization using this approach is denoted by Trankit-WT and Trankit-ST, respectively.

Trankit uses the detected words/tokens and sentences for POS tagging at sentence-level. For a given input sentence s , the transformers-generated representation of each of its words is computed by aggregating the representation of its word pieces. These distributed representational vectors

^{ee}<https://trankit.readthedocs.io/en/latest/> - Last visited: 21-October-2021.

are further fed to a *softmax* layer to predict the most probable POS tag for each word in the input sentence *s*. POS tagging using this approach is denoted by using Trankit-POST.

5. Proposed data set for UNLT

5.1 Data set for Urdu word tokenization

Testing data: Another key element of our research is to develop a large benchmark data set, for the evaluation of our proposed UNLT-WT approach (see Section 4.1.2). The process of developing a benchmark test data set is divided into three steps: (i) raw text collection, (ii) cleaning, and (iii) annotation.

In the first phase, raw data are collected from various online sources (BBC Urdu,^{ff} Express news,^{gg} Urdu Library,^{hh} Urdu Point,ⁱⁱ Minhaj Library,^{jj} Awaz-e-Dost^{kk} and Wikipedia^{ll}) by using a Web crawler.^{mm} The collected raw data are free and publicly available for research purposes and belongs to following genres: Commerce, Entertainment, Health, Weather, Science and Technology, Sports, Politics and Religion. This collected text contains 61,152 tokens.

In the next phase of the test data set creation process, the collected raw text was preprocessed (see Section 5.4), which resulted in the removal of 2152 tokens. The remaining cleaned data are composed of 59,000 tokens (3583 sentences).

The quality of evaluation of an Urdu word tokenization approach depends on the annotation quality of the test data set because inconsistent and noisy annotations deteriorate the model's performance. Thus, the annotations were performed by three different annotators (D, E, and F). All the annotators are native speakers of Urdu. The annotation process is further divided into three phases: (i) training, (ii) annotation, and (iii) inter-rater agreement calculation and conflict resolution.

In the training phase, two annotators (D and E) annotated a subset of 58 sentences. After that, the interannotator agreement was computed for these sentences, and conflicting tokens were discussed to further improve the annotation quality. In the annotation phase, the remaining test data set comprising of 3525 sentences was annotated by annotators D and E. After the annotation phase, the inter-rater reliability score was computed for the entire test data set of 59,000 tokens. We obtained the interannotator agreement of 86.3% as the annotators had agreement on 50,917 pairs. The Kappa coefficient was computed to be 78.09%, which is considered as good, considering the levels of difficulty for classifying the merge (space omission) and compound words (space insertion) into single or multiple tokens (see Section 3.1). Furthermore, the conflicting tokens were annotated, and decisions resolved by the third annotator F, which resulted in a gold standard UNLT-Word Tokenizer-Test (UNLT-WT-Test) data set.

Table 3 shows the type-token ratio of the UNLT-WT-Test data set that have a total of 59,000 tokens and 5849 types. The UNLT-WT-Test data set is stored in the standard "txt" format and is free and publicly available for research purposes (for license and URL see Section 7).

Training data: The training data set for our proposed Urdu word tokenizer was created by using a subset of the HC Corpus (Christensen 2014). To develop a gold standard training data set, two

^{ff}<http://www.bbc.com/urdu> terms of use: <https://www.bbc.com/urdu/institutional-37588278> - Last checked: 29-June-2019.

^{gg}<http://www.express.pk/> - Last checked: 28-October-2016.

^{hh}<http://www.urdulibrary.org/> - Last checked: 02-November-2016.

ⁱⁱ<http://www.urduweb.org/planet/> - Last checked: 08-November-2016.

^{jj}<http://www.minhajbooks.com/urdu/control/> - Last checked: 08-November-2016.

^{kk}<http://awaz-e-dost.blogspot.co.uk/> - Last checked: 08-November-2016.

^{ll}<https://ur.wikipedia.org/wiki/> - Last checked: 08-November-2016.

^{mm}<https://lindat.mff.cuni.cz/repository/xmlui/handle/11858/00-097C-0000-0023-65A9-5> - Last Checked: 18-December-2016.

Table 3. Domain wise statistics of the UNLT-WT-Test data set

Domain	Tokens	Types	Domain	Tokens	Types
Commerce	7254	663	Sports	6868	691
Entertainment	8578	937	Politics	9627	777
Health	6765	651	Religion	5553	556
Weather	6606	756	S&T ^a	7749	823

^aScience and Technology

million tokens were randomly selected from the following domains: Politics, Culture, Crime & Law, Fashion, Religion, Business & Economy, Science & Technology, Sports, Weather, Education, Health, Entertainment.

After preprocessing (see Section 5.4) the collected raw data, the resulting data set contained 1.65 million tokens. The preprocessed text is used to create the gold standard training data set. In the first step, the text is tokenized on the basis of space. After that, a human annotator manually corrected the improperly tokenized words generated in the first step. The final benchmark training data set (hereafter called UNLT-WT-Train data set) is composed of 1.65 million tokens.

The UNLT-WT-Train data set was used to generate N -grams using the approach described in Jurafsky and Martin (2014). Furthermore, we count the occurrences of each unique N -gram type, resulting in a total of 1,335,263 N -gram pairs with the following statistics: *tri*-grams: 636,765, *bi*-grams: 494,988, and *uni*-grams: 203,510.

5.2 Data set for Urdu sentence tokenization

For the evaluation of our proposed Urdu sentence tokenizer, we created a benchmark data set (hereafter called UNLT-ST data set) by following three steps: (i) raw Urdu text collection, (ii) preprocessing of raw data, and (iii) annotation. To construct the UNLT-ST data set, in the first step, we used a Web crawlerⁿⁿ to extract raw Urdu text of 12.5K sentences from online sources (see Section 5.1) including BBC Urdu, Express news, Urdu Library, Urdu Point, Minhaj Library, Awaz-e-Dost, and Wikipedia. These sources allow their text (content) to be freely used for research purposes. To make the data set more realistic, we extracted the raw text of different domains and genres including Sports, Politics, Blogs, Education, Literature, Entertainment, Science, Religion, Fashion, Weather, Entertainment, Fiction, Health, Law, and Business. BBC Urdu is our largest source of text collection, which contains 3358 sentences, while the Urdu Point is the smallest one containing 1157 sentences. Statistics of sentences collected from other sources are Awaze-e-Dost: 1457, Express news: 1557, Minhaj library: 1657, Urdu library: 1357, and Wikipedia: 1957 sentences.

In the second step, the raw data were preprocessed (see Section 5.4), which resulted in the removal of 2500 sentences. The remainder of the 10,000 clean sentences are distributed as follows: Awaz-e-Dost: 1200, BBC Urdu: 2606, and Express News: 1297, Minhaj Library: 1303, Urdu Library: 1119, Urdu Point: 948, and Wikipedia: 1527 sentences.

In the third step, the preprocessed text containing 10,000 cleaned sentences was manually tokenized by three annotators (G, H, and I). All the annotators are native speakers of Urdu and have good knowledge about Urdu sentence tokenization task. Furthermore, the annotation process was split into three phases: (i) training, (ii) annotation, and (iii) inter-rater agreement and conflict resolution.

ⁿⁿ<https://lindat.mff.cuni.cz/repository/xmlui/handle/11858/00-097C-0000-0023-65A9-5> - Last checked: 19-October-2016.

During the training phase, two annotators (G and H) annotated 1500 sentences. Subsequently, the interannotator agreement was computed for these sentences, and conflicting sentences were discussed to further improve the annotation quality. Further, during the annotation phase, the remaining 8500 sentences were manually annotated by annotators (G and H). In the third phase, the inter-rater agreement score was computed for all 10,000 sentences. We achieved an inter-rater agreement of 89.34%, as the annotators agreed upon 8934 sentences. Moreover, the Kappa coefficient was computed to be 81.83% (Cohen 1968). The conflicting 1066 sentences were annotated by the third annotator (I) for conflict resolution, and this judgement was considered as decisive, resulting in the gold standard UNLT-ST Training/Test data set.

The UNLT-ST Training/Test data set consists of 10,000 sentences. In our proposed test data set, 6469 period markers are SBM, while 536 are NSBM; 531 exclamation marks are SBM and 198 are NSBM; 421 question marks are SBM and 17 are NSBM; 197 double quotes, are SBM and 9 are NSBM; 382 SBM markers are #, @, \$, *; the remaining 2000 sentence are without any SBM. As can be noted from these statistics, our proposed UNLT-ST Training/Test data set contains both SBM and NSBM for different characters as well as sentences without any SBM, which makes the data set much more realistic and challenging. The UNLT-ST Training/Test data set is saved in standard “txt” format (for licensing and URL see Section 7)

5.3 Data set for Urdu POS tagging

This section describes the creation of a large data set (hereafter called UNLT-POS data set) for the training and testing of the Urdu POS taggers. The data set creation process was accomplished in three steps: (i) raw text collection, (ii) cleaning process, and (iii) annotation process.

To construct a gold-standard Urdu POS tagging data set, in the first step, a Web crawler (see Section 5.1) was used to extract Urdu text of 239,834 words (14,137 sentences) from various online sources (see Section 5.1) including BBC Urdu, Express news, Urdu library, Urdu point, Minhaj library, Awaz-e-Dost, and Wikipedia. To make the data set more realistic, the raw data are from various domains: Sports (23,153), Politics (33,944), Blogs (10,976), Education (12,845), Literature (9045), Entertainment (13,946), Science and Technology (17,683), Fashion (10,463), Weather (9459), Business (17,328), Commerce (10,496), Showbiz (19,503), Fictions (8678), Health (12,783), Law (8185), and Religion (21,347).

The raw data were preprocessed (see Section 5.4), which resulted in 200,000 words. The domain and genre distribution of these words is Sports (20,128), Politics (26,145), Blogs (9428), Education (10,742), Literature (8756), Entertainment (10,560), Science and Technology (13,143), Fashion (9758), Weather (8996), Business (14,418), Commerce (9710), Showbiz (16,228), Fictions (8084), Health (11,584), Law (6952), and Religion (15,368).

The UNLT-POS data set was created using a manual approach. In the first step, a total of 2000 tokens were POS tagged using the CLE online POS tagger^{oo} to train annotators. Manual inspection of the tagged data showed that a reasonable number of words are incorrectly tagged, particularly proper nouns, common nouns, verbs, auxiliaries, pronouns, adjectives, cardinal nominal modifiers, adverbs, conjunctions, participles, interjections, and foreign fragment. In the second training step, three annotators (A, B, and C) manually annotated^{pp} the automatically tagged data. Annotators A and B initially annotated 2000 tokens. An interannotator agreement was calculated for these tokens, and conflicting tagged tokens were discussed to further improve the annotation quality. After the training phase, the 200,000 words was manually annotated by annotators A and B, and the interannotator agreement was computed on the entire data set.

^{oo}<http://182.180.102.251:8080/tag/> - Last checked: 06-August-2016.

^{pp}In the training annotation process, the tag assigned by the CLE online POS tagger is retained if the annotator determines that it is correct; otherwise, the annotator replaces it with the correct POS tag.

Table 4. Statistics of three different training/testing data sets for evaluating the performance of Urdu POS taggers

Dataset		Training set	Testing set
UNLT-POS	Tokens	180,000	20,000
	Types	16,742	2124
	Unknown tokens	–	1948
	Unknown types	–	246
UNLT-POS-Moderate	Tokens	120,000	20,000
	Types	14,843	2457
	Unknown tokens	–	2078
	Unknown types	–	273
UNLT-POS-Small	Tokens	60,000	20,000
	Types	9538	2801
	Unknown tokens	–	3024
	Unknown types	–	311

An interannotator agreement of 85.7% was obtained. The Kappa coefficient was computed to be 77.41% (Cohen 1968). The conflicting tokens were annotated by the third annotator, resulting in a gold-standard UNLT-POS training/testing data set saved in “txt” format. As far as we are aware, our UNLT-POS training/testing data set is the largest manually POS tagged Urdu data set, free and publicly available for research purposes (for license and URL see Section 7).

For experiments presented in this study, the UNLT-POST gold-standard data set is randomly divided into two different data sets: (i) consisting of 60K training and 20K of test data (henceforth UNLT-POS-Small training/testing data set, respectively) and (ii) consisting of 120K training and 20K for testing (henceforth UNLT-POS-Moderate training/testing data set, respectively).

The detailed statistics of different train/test data sets are shown in Table 4. The rows “Unknown Tokens” and “Unknown Types” of Table 4 represent the count of total tokens and types (unique tokens), respectively, not seen in the different UNLT-POS training/testing data sets. It has been observed that each test data set holds 9% to 11% words that are unknown with respect to the training data. These figures are a little higher as compared to the several European languages (Evangelos and George 1995). However, Table 5 shows the detailed statistics of most frequent POS tags of the UNLT-POS testing data set.

5.4 Preprocessing

In this study, various data sets have been used, all these data sets (see Sections 5.1, 5.2, 5.3) are preprocessed as follows. Text in a data set is cleaned by removing multiple spaces, duplicated text, diacritics as they are optional (only used for altering pronunciation Mukund *et al.* 2010), and HTML tags. Moreover, noise from the data is removed by discarding ASCII and invalid UTF-8 characters, emoticons, asterisks, bullets, right and left arrows (Jawaid *et al.* 2014). Further, only

Table 5. Statistics of most frequent POS tags of UNLT-POS testing data set

POS Tag ^a	TC ^b	UT ^c	POS tag	TC	UT
NN	1,764	123	AUXA	1,023	0
PSP	1,572	0	NNP	1,243	398
VBF	1,129	192	RB	826	63
JJ	1,315	91	AUXT	639	3

^aNN, Common Noun; PSP, Postposition; VBF, Main Verb Finite; JJ, Adjective; AUXA, Aspectual Auxiliary; NNP, Proper Noun; RB, Common Adverb; AUXT, Tense Auxiliary.

^bToken Count.

^cUnknown tokens.

Table 6. Results obtained on UNLT-WT-Test data set using various techniques

Technique	Precision	Recall	F_1	Accuracy
UNLT-WT-SP	0.55	0.52	0.54	0.61
UNLT-WT	0.91	0.87	0.89	0.92
Durani's	0.18	0.20	0.19	0.49
CLE's	0.58	0.56	0.57	0.73
Trankit-WT	0.73	0.69	0.71	0.80

sentences with three or more words were kept.⁹⁹ A language detection tool^{rr} is used to discard foreign words and a text normalization tool.^{ss}

6. Results and analysis of proposed UNLT tools

6.1 Results of Word tokenizer

Table 6 presents precision, recall, F_1 , and accuracy results when training on UNLT-WT-Train data set and testing on the UNLT-WT-Test (see Section 5.1) for Urdu word tokenization task by using various approaches. UNLT-WT-SP refers to results obtained using the space-based tokenization approach. UNLT-WT refers to results obtained using our proposed approach (see Section 4.1.2) for Urdu word tokenization. Durani's word tokenizer (see Section 2.1). Whereas CLE's word tokenizer refers to an online tokenizer.^{tt} The Trankit-WT is a deep learning model that implements a light-weight XLM-Roberta transformer (see Section 4.4).

Overall, the best results are obtained by using our proposed UNLT-WT approach (precision = 0.91, recall = 0.87, F_1 = 0.89, and accuracy = 0.92). These results show that UNLT-WT is the most appropriate method for Urdu word tokenization on the UNLT-WT-Test data set. Furthermore, this also shows that combining maximum matching, dictionary lookup, and statistical N -gram MLE along with smoothing estimation are helpful in getting good performance on UNLT-WT-Test data set for Urdu word tokenization task. However, the highest F_1 score of

⁹⁹This is calculated by dividing the total words in data set by the total number of sentence disambiguation markers.

^{rr}<https://lindat.mff.cuni.cz/repository/xmlui/handle/11858/00-097C-0000-0023-65A9-5> - Last checked: 15-October-2016.

^{ss}Text normalization tool can be downloaded from <http://www.cle.org.pk/software/langproc/urduNormalization.htm> - Last checked: 15-October-2016 is used to keep the Unicode of the characters consistent.

^{tt}<http://www.cle.org.pk/clestore/segmentation.htm> - Last visited: 18-Dec-2019. Tokenized up to 100 words at one time and implementation details are not available. The online link refers three papers but does not describe which one of them is used for the creation of CLE Urdu word tokenizer.

Table 7. Results obtained by using various sentence tokenization approaches on UNLT-ST-Train/Test data set

Technique	Precision	Recall	F_1	Error rate
UNLT-ST-PQEQM	0.89	0.19	0.31	0.83
UNLT-ST-RB	0.81	0.74	0.77	0.29
UNLT-ST-ML	0.90	0.92	0.91	0.09
Trankit-ST	0.86	0.83	0.85	0.14

0.89 for the word tokenization task indicates that Urdu word tokenization is a challenging task leaving room for further improvement.

As expected, the overall results are lower for the baseline space-based tokenization UNLT-WT-SP (precision = 0.55, recall = 0.52, F_1 = 0.54, and accuracy = 0.61) and approach, on UNLT-WT-Test data set. Durani's word tokenizer reports an accuracy of 0.49, precision of 0.18, recall of 0.20, and F_1 = 0.19. Furthermore, the CLE's Urdu word tokenizer has shown precision = 0.58, recall = 0.56, F_1 = 0.57, and accuracy = 0.73. Whereas the Trankit-WT has produced the following results: precision = 0.73, recall = 0.69, F_1 = 0.71, and accuracy = 0.80. This highlights the fact that the UNLT-WT-SP, Durani's, CLE's approaches are not appropriate for Urdu word tokenization tasks.

The comparison of our proposed methods significantly outperformed the recently reported deep neural network-based word-piece based tokenization (Trankit-WT). There could be various possible reasons for this notable difference in the scores including simple space-based splitting of the raw input text as a preprocessing step in Trankit-WT seems to be a major contributor in producing low scores for word tokenization. This preprocess space-splitting could be efficient for many languages like English, Spanish, etc. but not suitable for the Urdu language. As has been described previously that spaces play a major role in word tokenization (see Section 3.1). We have further observed that Trankit-W is very poor in capturing multiwords.

While analyzing the errors of the proposed UNLT-WT approach, we observed that it does not explicitly handle unknown words for space omission, and this resulted in splitting an unknown Urdu morpheme into smaller morphemes. For instance, the word کثیراللسان (KSYR ALLSAN, "multilingual") erroneously split into کثی (KSY), را (RA), لیس (LLS), and ان (AN). Likewise, it might be less appropriate when a word is a combination of known and unknown morphemes, for instance, شہبازکو جانیدو (SHBAZ KO JANE DO, "let the Shahbaz go"). For space insertion, some compound words were not present in compound words dictionary, another major cause of incorrect word tokenization.

6.2 Results of sentence tokenizer

Table 7 presents precision, recall, F_1 , and error rate results^{uu} on the UNLT-ST Train/Test data set^{vv} (see Section 5.2) for various Urdu sentence tokenization approaches^{ww} (for other two, see Section 4.2.2). Trankit-ST is a light-weight XLM-Roberta transformer for sentence tokenization (see Section 4.4).

Overall, the best results are obtained using our proposed UNLT-ST-ML approach, precision = 0.90, recall = 0.92, F_1 = 0.91, and error rate = 0.09. The lowest results are obtained using a baseline UNLT-ST-PQEQM (precision = 0.89, recall = 0.19, F_1 = 0.31, and error rate = 0.83).

^{uu}For ML approach, we have used 10-fold cross-validation.

^{vv}For baseline and rule-based approaches, we have used the entire UNLT-ST Train/Test data set as a test data set.

^{ww}Baseline UNLT-ST-PQEQM—tokenization on the basis of "period," "question mark," "exclamation mark," and "double quotes" characters.

Whereas another proposed rule-based UNLT-ST-RB approach shows a precision of 0.81, recall of 0.74, F_1 of 0.77, and error rate of 0.29. This shows that combining various features in SVM is helpful in producing a good performance on the UNLT-ST Train/Test data set. Trankit-ST a deep learning-based approach has produced results as follows: precision = 0.86, recall = 0.83, F_1 = 0.85, and error rate = 0.14. It is worth mentioning here that Trankit-S has produced more accurate results than two other sentence tokenization approaches, that is, UNLT-ST-RB/PQEQM. However, the highest F_1 score of 0.91 for sentence tokenization task indicates that Urdu sentence tokenization is a challenging task, and there is still room for further improvement.

For the UNLT-ST-PQEQM approach, which uses different characters as sentence boundary indicators, the performance is high in term of precision of 0.89. The likely reason for this is the majority of sentences in Urdu text are terminated using these characters. However, other evaluation measures show very low results. This highlights the fact that these characters alone are not suitable for the Urdu sentence tokenization task. As far as the rule base (UNLT-ST-RB) approach is concerned, it shows good results as compared to the baseline approach. This shows that combining various heuristics, regular expressions, and dictionary lookup is helpful in producing a good performance on the UNLT-ST Test data set. However, this method also fails to detect sentence boundary where sentences are ended without any SBM. While manually analyzing the errors of the proposed UNLT-ST-RB approach, we came across some scenarios where our proposed approach failed to accurately tokenize sentences. It was found that NSBM including “:”, “||”, “\$”, “*”, “@”, and “#” are the major reasons for incorrect tokenization of sentences. Moreover, the period used between different abbreviations also caused misclassification. Finally, all those sentences that end without any sentence boundary marker are also declassified.

6.3 Results of POS tagger

Table 8 presents accuracy results when trained and tested on the UNLT-POS-Small, UNLT-POS-Moderate, UNLT-POS test data sets (see Section 5.3) for the Urdu POS tagging tasks by using different models (see Sections 4.3.2 and 4.4).

Overall best results are obtained using our proposed T-HMM-GT-Suf-MA followed by T-HMM-GT-Suf-MA and MEn-Suf-MA POS tagging models, 95.59%, 95.14%, and 94.20%, respectively. This shows that combining various stochastic and smoothing techniques with language-dependent features is helpful in producing a very good performance on the UNLT-POS test data set. The highest accuracy score of 95.59% indicates that the Urdu POS tagging task is challenging, and there is still room for improvement. It can also be noted from these results that our proposed POS tagging approach (T-HMM-GT-Suf-MA) outperforms both baseline approaches BL-MFT^{xx} (accuracy = 84.72%) and BL-CLE^{yy} (Ahmed *et al.* 2014) (which uses Decision Trees along with a smoothing technique of Class Equivalence) (accuracy = 88.45%) on UNLT-POS test data set. Furthermore, Trankit-POST has produced an accuracy of 92.67%. This has good accuracy as compared to baseline and several other approaches. The reason for using the BL-CLE model as a baseline approach is that, currently, this is the only POS tagger available for Urdu which uses CLE Urdu POS tagset (see Section 4.3.1). Therefore, we can compare the results of CLE Urdu POS tagger with our proposed UNLT-POS tagger.

We can further observe that the *tri*-gram HMM-based models can produce good results if incorporated with linear interpolation, suffix, as well as Morphological Information (MI). Certainly, using MI along with linear interpolation gives better results as compared to suffix, but what is significant to note, using all the information together improved the accuracy of the models, T-HMM-LI-Suf-MA: 95.14%, T-HMM-LaE-Suf-MA: 93.74%, T-HMM-LiE-Suf-MA: 93.97%, T-HMM-GT-Suf-MA: 94.48%, T-HMM-KN-Suf-MA: 95.59%, and MEn-Suf-MA:

^{xx}In this, each word in the test data will be assigned the POS tag based on the most frequent POS tag in the training data.

^{yy}<http://182.180.102.251:8080/tag/> - Last checked: 17-July-2018.

Table 8. Results obtained using various POS tagging models based on several approaches on different POS test data sets

Approaches ^a	Model	Accuracy ^b		
		D1	D2	D3
Most frequent tag	BL-MFT	–	–	84.72
Decision Tree	BL-CLE	–	–	88.45
Trankit-POST	T-DL	–	–	92.67
<i>tri</i> -gram HMM, LI	T-HMM-LI	67.14	80.34	87.34
<i>tri</i> -gram HMM, LI, suffix	T-HMM-LI-Suf	83.23	87.91	91.53
<i>tri</i> -gram HMM, LI, MI	T-HMM-LI-MA	88.37	90.39	92.27
<i>tri</i> -gram HMM, LI, suffix, MI	T-HMM-LI-Suf-MA	90.87	93.76	95.14
<i>tri</i> -gram HMM, LaE	T-HMM-LaE	65.97	79.14	85.92
<i>tri</i> -gram HMM, LaE, suffix	T-HMM-LaE-Suf	80.42	86.39	89.98
<i>tri</i> -gram HMM, LaE, MI	T-HMM-LaE-MA	87.88	89.74	90.19
<i>tri</i> -gram HMM, LaE, suffix, MI	T-HMM-LaE-Suf-MA	89.04	91.64	93.74
<i>tri</i> -gram HMM, LiE	T-HMM-LiE	66.98	80.02	86.89
<i>tri</i> -gram HMM, LiE, suffix	T-HMM-LiE-Suf	82.78	87.13	90.93
<i>tri</i> -gram HMM, LiE, MI	T-HMM-LiE-MA	88.13	90.02	91.69
<i>tri</i> -gram HMM, LiE, suffix, MI	T-HMM-LiE-Suf-MA	90.23	92.59	93.97
<i>tri</i> -gram HMM, GT	T-HMM-GT	67.02	80.10	87.01
<i>tri</i> -gram HMM, GT, suffix	T-HMM-GT-Suf	82.98	87.27	91.05
<i>tri</i> -gram HMM, GT, MI	T-HMM-GT-MA	88.19	90.07	91.90
<i>tri</i> -gram HMM, GT, suffix, MI	T-HMM-GT-Suf-MA	90.57	93.03	94.48
<i>tri</i> -gram HMM, KN	T-HMM-LiE	67.24	80.36	87.51
<i>tri</i> -gram HMM, KN, suffix	T-HMM-LiE-Suf	83.23	87.98	91.70
<i>tri</i> -gram HMM, KN, MI	T-HMM-LiE-MA	88.42	90.52	92.39
<i>tri</i> -gram HMM, KN, suffix, MI	T-HMM-LiE-Suf-MA	91.02	93.99	95.59
MaEn, CW, WN	MEn	80.59	84.92	88.31
MaEn, CW, WN, suffix	MEn-Suf	84.43	88.06	92.56
MaEn, CW, WN, MI	MEn-MA	88.32	89.49	93.11
MaEn, CW, WN, suffix, MI	MEn-Suf-MA	90.26	93.31	94.20

^a LI: Linear Interpolation, MI: Morphological Information, LaE: Laplace Estimation, LiE: Lidstone's Estimation, GT: Good-Turing, KN-Kneser-Ney MaEn: Maximum Entropy, CW: Context Window, WN: Word Number

^b D1: UNLT-POS-Small training/testing data set, D2: UNLT-POS-Moderate training/testing data set, D3: UNLT-POS training/testing data set

94.20%. Furthermore, it can be observed, T-HMM-LI, T-HMM-LaE, T-HMM-LiE, T-HMM-GT, and T-HMM-KN produce accuracies of 87.34%, 85.92%, 86.89%, 87.01%, and 87.51%, respectively, on the UNLT-POS data set. For the case of MEn, the reported accuracy is 88.31%. One important observation here is that by using smoothing and language-dependent features, the proposed Urdu POS tagging accuracies can be improved as compared to BL-MFT and BL-CLE models.

It can also be observed from the Table that T-HMM-GT performs better than the other four models T-HMM-LI, T-HMM-LaE, T-HMM-LiE, T-HMM-GT, and T-HMM-KN on UNLT-POS, UNLT-POS-Small, and UNLT-POS-Moderate test data sets. Moreover, the accuracy of T-HMM-LaE model is slightly poorer than the other HMM-based models (T-HMM-LI, T-HMM-LiE, T-HMM-GT, and T-HMM-KN), with UNLT-POS-Small data due to model overfitting. However, such discrepancies are alleviated with the increase of training data (UNLT-POS-Moderate and UNLT-POS training data sets).

It has been further observed that language-dependent features increased the accuracy of the models to a certain extent, even if trained on a UNLT-POS-Moderate training data set. However, with different features along with smoothing, the increase in the model accuracy is higher when training data are smaller. For instance, T-HMM-LI-MA and T-HMM-LI-Suf models improved around 16%, 7%, and 4% and 21%, 10%, and 5%, respectively, over the T-HMM-LI models, for UNLT-POS, UNLT-POS-Small, and UNLT-POS-Moderate test data sets.

From the above observations, it can be concluded that using MI and suffix increases in the model accuracy are higher for UNLT-POS-Small and UNLT-POS-Moderate training data sets. It is also important to note that the T-HMM-KN-MA models give an approximate improvement of around 5%, 7%, and 1% over the T-HMM-KN-Suf model for UNLT-POS-Small, UNLT-POS-Moderate, and UNLT-POS training data set, respectively. However, integrating all of them, an improvement has been observed in T-HMM-KN-Suf-MA models which are 8%, 6%, and 3% improved with respect to T-HMM-KN-Suf model in case of UNLT-POS-Small, UNLT-POS-Moderate, and UNLT-POS training data set. It can also be noticed that similar results have been observed for the other two (T-HMM-LaE, T-HMM-LiE, T-HMM-GT, and T-HMM-LI) HMM-based models. However, T-HMM-LiE performed better than the T-HMM-LaE model, but with the higher training data, the performance of these models is somewhat comparable.

MEn models outperform all others with smaller training data, but contrasting results have been observed with large training data. It is worth noting that MEn along with suffix and morphological information has positive effects with poor resources. Our results show the T-HMM-GT-Suf-MA and MEn-Suf-MA are more accurate than others, providing support for our further analysis based on such models.

Table 9 shows cases where the MEn-Suf-MA model performs better than T-HMM-KN-Suf-MA, by comparing the accuracies of open class tags for known and unknown words on the UNLT-POS testing data set. Our results show that the T-HMM-KN-Suf-MA model shows poor accuracy while predicting proper nouns (NNP) over the MEn-Suf-MA model. Mostly, the proper nouns (NNP) in T-HMM-KN-Suf-MA model are erroneously classified as an adjective (JJ). Furthermore, it is worth noting again that in Urdu, there is no discrimination between upper and lower-case characters, also using an adjective as a proper noun is frequent in Urdu, for example, کبر (KBYR, “big”) and صغر (SGHYR, “small”). Another reason for misclassification in tagging of the proper nouns is that many of them end with negation marker or pronoun, for example, the ناگنہ (“Nagyna”) end with the نہ (NH, “no”) or the NNP نازہ (“Nazyh”) which end with the یہ (YH, “this”), a pronoun. These errors need further investigation.

Similarly, in the case of common adverb (RB), the accuracy of the MEn-Suf-MA model is approximately 2% and 6% higher for known and unknown words, respectively. However, the performance of MEn-Suf-MA for all other open class tags did not improve over the T-HMM-KN-Suf-MA. It is further observed that with the increase of unknown words, the accuracy of the MEn model has reduced.

Table 9. Accuracies of open class tags on UNLT-POS testing data set using T-HMM-KN-Suf-MA and MEn-Suf-MA

Tag	T-HMM-KN-Suf-MA		MEn-Suf-MA	
	Known	Unknown	Known	Unknown
NN	95.10	80.34	92.32	78.23
NNP	73.98	56.09	76.56	70.74
JJ	92.02	63.58	89.54	61.97
RB	81.98	57.98	84.45	64.33
VBF	93.14	73.01	92.47	72.03

The most prominent causes for Urdu POS tag misclassification are that there is no clear distinction between noun and proper noun, dropping of words is also frequent, if a noun in a noun phrase is dropped the adjective becomes a noun in that phrase, the highly inflected nature of Urdu, and ambiguity between noun and verb is due to verbal nouns.

7. Conclusion and future directions

In this paper, we have described a novel Urdu Natural Language Toolkit by integrating word and sentence tokenizers as well as a POS tagger. Words are tokenized by coupling a rule-based morpheme matching method with a *tri*-gram stochastic language model, backed-off to *bi*-gram Maximum Likelihood Estimation supplemented by smoothing technique for unseen words. We also developed large compound word and morpheme dictionaries, which were used in our proposed Urdu word tokenizer. A large benchmark training and testing data sets are also created. The training data set comprises 1,361,179 *N*-grams (1.65 million tokens), whereas test data set contains 59,000 manually tokenized words. Moreover, we have compared our results with state-of-the-art deep learning methods for comparison purposes. The results show that our proposed Urdu word tokenizer obtained precision of 0.96, recall of 0.92, F_1 of 0.94, and accuracy of 0.97. The sentence tokenization approach is formed by using rules extracted from a large raw Urdu corpus, regular expressions, and dictionary lookup. Our proposed Urdu sentence tokenizer obtained promising results on a large and new generated gold-standard test data set (precision = 91.08%, recall = 94.14%, F_1 = 92.59%, and error rate = 6.85%). We have described 24 different stochastic models (and two baseline models) for the Urdu POS tagging task. Each of these models has a unique stochastic scheme supplemented with various language features and smoothing estimations. In addition, a large gold-standard training/testing data set was generated. Results show that the best performance is achieved by the T-HMM-KN-Suf-MA POS tagger which is a combination of *tri*-gram HMM, Kneser–Ney, suffix, and morphological information, achieving an accuracy of 95.59%. These resources have been made publicly available to the research community at <https://github.com/UCREL/UNLT> and <https://doi.org/10.17635/lancaster/researchdata/494>, under the terms of the Creative Commons Attribution 4.0 International License^{zz} and GNU General Public License v3.0.^{aaa}

Our next primary target for follow-up research work will be to extend the toolkit by developing and integrating various NLP tools. As far as the word tokenization is concerned, we aim to adopt machine learning approaches (conditional random field, maximum entropy, etc.) to learn

^{zz}<https://creativecommons.org/licenses/by/4.0/>- Last checked: 15-September-2020.

^{aaa}<https://www.gnu.org/licenses/gpl-3.0.en.html>- Last checked: 11-November-2021.

the morphological pattern of the valid morphemes (instead of morphemes lookup) and to handle out-of-vocabulary words in morpheme matching process of space omission problem. For the sentence tokenizer, we plan to develop a hybrid approach for Urdu sentence tokenization task, that is, rule-based along with artificial neural network. Finally, in the future, we aim to deal with unknown words and will adopt further statistical methods (CRF, deep learning, etc.) along with heuristic rules to increase the POS tagging accuracy for Urdu text.

Acknowledgements. This work has been supported by the COMSATS University Islamabad, Lahore Campus, Pakistan and Lancaster University, U.K. under the Split-Site Ph.D. programme.

References

- Abdelhamid A.A., Abdulla W.H. and MacDonald B. (2012). WFST-based large vocabulary continuous speech decoder for service robots. In *Proceedings of the International Conference on Imaging and Signal Processing for Healthcare and Technology (ISPHT'12)*, Baltimore, USA, pp. 150–154.
- Ahmed T., Urooj S., Hussain S., Mustafa A., Parveen R., Adeeba F., Hautli A. and Butt M. (2014). The CLE Urdu POS tagset. In *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland, pp. 2920–2925.
- Akita Y., Saikou M., Nanjo H. and Kawahara T. (2006). Sentence boundary detection of spontaneous Japanese using statistical language model and support vector machines. In *Proceedings of the International Conference on Spoken Language Processing (ICSLP'06)*, Pennsylvania, USA, pp. 1033–1036.
- Albared M., Omar N., Aziz M.J. and Nazri M.Z. (2010). Automatic part of speech tagging for Arabic: an experiment using bigram hidden Markov model. In *International Conference on Rough Sets and Knowledge Technology, (RSKT'10)*, Beijing, China, pp. 361–370.
- Anwar W., Wang X., Li L. and Wang X.L. (2007a). A statistical based part of speech tagger for Urdu language. In *IEEE International Conference on Machine Learning and Cybernetics (ICMLC'07)*, Hong Kong, China, vol. 6, pp. 3418–3424.
- Anwar W., Wang X., Li L. and Wang X.L. (2007b). Hidden Markov model based part of speech tagger for Urdu. *Information Technology Journal* 6(8), 1190–1198.
- Azimzadeh A., Arab M.M. and Quchani S.R. (2008). Persian part of speech tagger based on Hidden Markov Model. In *Proceedings of the 9th International Conference on the Statistical Analysis of Textual Data (JADT'08)*, Lyon, France, pp. 121–128.
- Bhat R.A. and Sharma D.M. (2012). A dependency treebank of Urdu and its evaluation. In *Proceedings of the Sixth Linguistic Annotation Workshop (LAW VI'12)*, Jeju, Republic of Korea, pp. 157–165.
- Bird S., Klein E. and Loper E. (2009). *Natural Language Processing with Python*. O'Reilly Media, Inc.
- Bird S., Klein E., Loper E. and Baldridge J. (2008). Multidisciplinary instruction with the natural language toolkit. In *Proceedings of the 3rd Workshop on Issues in Teaching Computational Linguistics (TeachCL'08)*, Ohio, USA, vol. 13, pp. 62–70.
- Bögel T., Butt M., Hautli A. and Sulger S. (2007). Developing a finite-state morphological analyzer for Urdu and Hindi. In *Proceedings of the Finite-State Methods and Natural Language Processing: 6th International Workshop (FSM/NLP'07)*, Potsdam, Germany, pp. 86–96.
- Brants T. (2000). TnT: a statistical part-of-speech tagger. In *Proceedings of the Sixth Applied Natural Language Processing Conference ANLP-2000*, Seattle, Washington, USA, pp. 224–231.
- Butt J.M. (1995). *The Structure of Complex Predicates in Urdu*. PhD Thesis, Center for the Study of Language (CSLI), department of linguistics, Stanford University.
- Chen S.F. and Goodman J. (1999). An empirical study of smoothing techniques for language modeling. *Computer Speech & Language* 13(4), 359–394.
- Christensen H. (2014). HC corpora. Available at <http://www.corpora.heliohost.org/> (accessed 5 March 2017).
- Christer S. (1996). Handling sparse data by successive abstraction. In *Proceedings of the 16th International Conference on Computational Linguistics (COLING'96)*, Copenhagen, Denmark, pp. 895–900.
- Cohen J. (1968). Weighted kappa: nominal scale agreement provision for scaled disagreement or partial credit. *Psychological bulletin* 70(4), 213–20.
- Conneau A., Khandelwal K., Goyal N., Chaudhary V., Wenzek G., Guzman F., Grave E., Ott M., Zettlemoyer L. and Stoyanov V. (2020). Unsupervised cross-lingual representation learning at scale. arXiv preprint arXiv:1911.02116, 1, 1–12.
- Cunningham H., Maynard D., Bontcheva K. and Tablan V. (2002). GATE: a framework and graphical development environment for robust NLP tools and applications. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, Philadelphia, USA, pp. 168–175.

- Cunningham H., Tablan V., Roberts A. and Bontcheva K. (2013). Getting more out of biomedical documents with GATE's full lifecycle open source text analytics. *PLoS Computational Biology* 9(2), e1002854.
- Curran J.R. and Clark S. (2003). Investigating GIS and smoothing for maximum entropy taggers. In *Proceedings of the 10th Conference on European Chapter of the Association for Computational Linguistics (EACL'03), Budapest, Hungary*, vol. 1, pp. 91–98.
- Dandapat S. (2008). Part of speech tagging and chunking with maximum entropy model. In *Proceedings of the IJCAI Workshop on Shallow Parsing for South Asian Languages (IJCAI'08), Hyderabad, India*, pp. 29–32.
- Daud A., Khan W. and Che D. (2016). Urdu language processing: a survey. *Artificial Intelligence Review* 47(3), 279–311.
- Dietzel A. and Maynard D. (2015). Climate change: a chance for political re-engagement? In *Proceedings of the Political Studies Association 65th Annual International Conference (PSA'15), Sheffield, UK*, pp. 1–19.
- Durrani N. and Hussain S. (2010). Urdu word segmentation. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the ACL, Los Angeles, California, USA*, pp. 528–536.
- Ekbal A., Haque R. and Bandyopadhyay S. (2008). Maximum entropy based Bengali part of speech tagging. *Advances in Natural Language Processing and Applications, Research in Computing Science (RCS) Journal* 33(8), 67–78.
- Evangolos D. and George K. (1995). Automatic stochastic tagging of natural language texts. *Computational Linguistics* 21(2), 137–63.
- Ferrucci D. and Lally A. (2004). UIMA: an architectural approach to unstructured information processing in the corporate research environment. *Natural Language Engineering* 10(3–4), 327–348.
- Fu G., Kit C. and Webster J.J. (2008). Chinese word segmentation as morpheme-based lexical chunking. *Information Sciences* 178(9), 2282–2296.
- Garside R. and Smith N. (1997). A hybrid grammatical tagger: CLAWS4. In *Corpus Annotation: Linguistic Information from Computer Text Corpora, Longman, London*, vol. 102(1), pp. 102–121.
- Giménez J. and Marquez L. (2004). SVMTool: a general POS tagger generator based on Support Vector Machines. In *Proceedings of the 4th International Conference on Language Resources and Evaluation (LREC'04), Lisbon, Portugal*, pp. 43–46.
- Gries S.T. and John N. (2014). *Research Methods in Linguistics*. Cambridge University Press, UK, pp. 257–287.
- Hardie A. (2003). Developing a tagset for automated part-of-speech tagging in Urdu. In Archer D., Rayson P., Wilson A. and McEnery T. (eds), *Proceedings of the Corpus Linguistics 2003 Conference. UCREL Technical Papers, Lancaster, UK*, vol. 16, pp. 298–307.
- Hardie A. (2004). *The Computational Analysis of Morphosyntactic Categories in Urdu*. PhD Thesis, Lancaster University, UK.
- Hautli A. and Sulger S. (2011). Extracting and classifying Urdu multiword expressions. In *Proceedings of the the ACL-HLT Student Session (ACL-HLT'11), Portland, OR, USA*, pp. 24–29.
- Hearst M.A., Dumais S.T., Osuna E., Platt J. and Scholkopf B. (1998). Support vector machines. *IEEE Intelligent Systems and their Applications* 13(4), 18–28.
- Javed I. (1985). *Nai Urdu Qawaid*. Urdu Development Board, New Delhi.
- Jawaid B., Kamran A. and Bojar O. (2014). A tagged corpus and a tagger for Urdu. In *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC'09), Reykjavík, Iceland*, pp. 2938–43.
- Jeffreys H. (1998). *The Theory of Probability*, 3rd Edn. Oxford University Press.
- Joshi N., Darbari H. and Mathur I. (2013). HMM based POS tagger for Hindi. In *Proceedings of the 2013 International Conference on Artificial Intelligence and Soft Computing (AISC'13), Bangalore, India*, pp. 341–349.
- Jurafsky D. and Martin J. (2014). *Speech & Language Processing*, 2nd Edn., vol. 3. Pearson, London.
- Khan S.A., Anwar W., Bajwa U.I. and Wang X. (2012). A light weight stemmer for Urdu language: a scarce resourced language. In *Proceedings of the 3rd Workshop on South and Southeast Asian Natural Language Processing (SANLP-COLING'12), Mumbai, India*, pp. 69–78.
- Kreuzthaler M. and Schulz S. (2015). Detection of sentence boundaries and abbreviations in clinical narratives. *BMC Medical Informatics and Decision Making* 15(2), 1–13.
- Kudo T. and Richardson J. (2018). SentencePiece: a simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (System Demonstrations), Brussels, Belgium*, pp. 66–71.
- Kudo T. (2018). Subword regularization: improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Long Papers), Melbourne, Australia*, pp. 66–75.
- Kwartler T. (2017). *Text Mining in Practice with R*. John Wiley & Sons, Ltd, Chichester, UK, pp. 237–269
- Lehal G.S. (2010). A word segmentation system for handling space omission problem in Urdu script. In *Proceedings of the 1st Workshop on South and Southeast Asian Natural Language Processing (WSSANLP), the 23rd International Conference on Computational Linguistics, Beijing, China*, pp. 43–50.
- Malik A. (2009). A hybrid model for Urdu Hindi translation. In *Proceedings of the Named Entities WorkShop (NEWS'09), Joint Conference of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, (ACL-IJCNLP'09) Singapore*, pp. 177–185.

- Manning C.D. and Schütze H.** (1999). *Foundations of Statistical Natural Language Processing*, vol. 999. MIT Press, Cambridge, MA.
- Manning C.D., Surdeanu M., Bauer J., Finkel J.R., Bethard S. and McClosky D.** (2014). The Stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations, Baltimore, Maryland, USA*, pp. 55–60.
- Maynard D., Greenwood M.A., Roberts I., Windsor G. and Bontcheva K.** (2015). Real-time social media analytics through semantic annotation and linked open data. In *Proceedings of the ACM Web Science Conference (WebSci'15), Oxford, UK*, pp. 46–48.
- Muaz A., Ali A. and Hussain S.** (2009). Analysis and development of Urdu POS tagged corpus. In *Proceedings of the 7th Workshop on Asian Language Resources (ALR'7), Suntec, Singapore*, pp. 24–29.
- Mukund S., Srihari R. and Peterson E.** (2010). An information-extraction system for Urdu—a resource-poor language. *ACM Transactions on Asian Language Information Processing (TALIP)* 9(4), 1–43.
- Naz F., Anwar W., Bajwa U.I. and Munir E.** (2012). Urdu part of speech tagging using transformation based error driven learning. *World Applied Sciences Journal (WASJ)* 16(3), 437–448.
- Nguyen M., Lai V., Veyseh A.P.B. and Nguyen T.H.** (2021). Trankit: a light-weight transformer-based toolkit for multi-lingual Natural Language Processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations, Online*, pp. 80–89.
- Platts J.T.** (1909). *A Grammar of the Hindustani or Urdu Language*. Crosby Lockwood and Son, London, republished in 2002 by Sang-e-Meel Publications, Lahore.
- Rabiner L.R.** (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE* 77(2), 257–285.
- Raj S., Rehman Z., Rauf S., Siddique R. and Anwar W.** (2015). An artificial neural network approach for sentence boundary disambiguation in Urdu. *The International Arab Journal of Information Technology* 12(4), 395–400.
- Rashid R. and Latif S.** (2012). A dictionary based Urdu word segmentation using maximum matching algorithm for space omission problem. In *Proceedings of the International Conference on Asian Language Processing (IALP'17), Hanoi, Vietnam*, pp. 101–104.
- Ratnaparkhi A.** (1996). A maximum entropy model for part-of-speech tagging. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP'96), New Jersey, USA*, vol. 1, pp. 133–142.
- Rehman Z., Anwar W. and Bajwa U.I.** (2011). Challenges in Urdu text tokenization and sentence boundary disambiguation. In *Proceedings of the 2nd Workshop on South Southeast Asian Natural Language Processing (WSSANLP'11), Chiang Mai, Thailand*, pp. 40–45.
- Rehman Z. and Anwar W.** (2012). A hybrid approach for Urdu sentence boundary disambiguation. *The International Arab Journal of Information Technology* 9 (3), 250–255.
- Rehman Z., Anwar W., Bajwa U.I., Xuan W. and Chaoying Z.** (2013). Morpheme matching based text tokenization for a scarce resourced language. *PLoS One* 8(8), e68178.
- Riaz K.** (2012). Comparison of Hindi and Urdu in computational context. *International Journal of Computational Linguistics and Natural Language Processing (IJCLNLP)* 1(3), 92–97.
- Riaz K.** (2010). Rule-based named entity recognition in Urdu. In *Proceedings of the 2010 Named Entities Workshop (NEWS'10), Uppsala, Sweden*, pp. 126–135.
- Rush A.M., Chopra S. and Weston J.** (2015). A neural attention model for abstractive sentence summarization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP'15), Lisbon, Portugal*, pp. 1–11.
- Saeed A., Nawab R.M.A., Stevenson M. and Rayson P.** (2018). A word sense disambiguation corpus for Urdu. *Language Resources and Evaluation (LRE)* 1(3), 1–22.
- Sajjad H.** (2007). *Statistical Part of Speech Tagger for Urdu*. Unpublished MS Thesis, National University of Computer and Emerging Sciences, Lahore, Pakistan.
- Sajjad H. and Schmid H.** (2009). Tagging Urdu text with parts of speech: a tagger comparison. In *Proceedings of the 12th Conference of the European Chapter of the ACL, Athens, Greece*, pp. 692–700.
- Schmid H.** (1994b). Probabilistic part-of-speech tagging using decision trees. In *Proceedings of the International Conference on New Methods in Language Processing (NeMLaP), Manchester, UK*, vol. 12, pp. 44–49.
- Schmid H. and Laws F.** (2008). Estimation of conditional probabilities with decision trees and an application to fine-grained POS tagging. In *Proceedings of the 22nd International Conference on Computational Linguistics (COLING'08), Manchester, UK*, vol. 1, pp. 777–784.
- Schmidt R.L.** (1999). *Urdu, an Essential Grammar (Routledge Essential Grammars)*, vol. 1. Psychology Press.
- Shafi J.** (2020). *An Urdu Semantic Tagger—Lexicons, Corpora, Methods, and Tools*. PhD Thesis, Lancaster University, UK.
- Sharjeel M., Nawab R.M.A. and Rayson P.** (2017). COUNTER: corpus of Urdu news text reuse. *Language Resources and Evaluation* 1(3), 777–803.
- Tafseer A.** (2009). Roman to Urdu transliteration using wordlist. In *Proceedings of the Conference on Language and Technology (CLT'09), Lahore, Pakistan*, pp. 1–8.

- Thede S.M. and Harper M.P.** (1999). A second-order hidden Markov model for part-of-speech tagging. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics on Computational Linguistics (ACL'99)*, College Park, Maryland, pp. 175–82.
- Vaswani A., Shazeer N., Parmar, Uszkoreit N., Jones L., Gomez A.N., Kaiser L. and Polosukhin I.** (2017). Attention Is All You Need. arXiv preprint arXiv:1706.03762 1, 1–5.
- Viterbi A.** (1967). Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory* **13**(2), 260–269.
- Wicaksono A.F. and Purwarianti A.** (2010). HMM based part-of-speech tagger for Bahasa Indonesia. In *Proceedings of the Fourth International MALINDO Workshop, Jakarta, Indonesia*, pp. 1–7.
- Yi C.** (2015). An English pos tagging approach based on maximum entropy. In *Proceedings of the International Conference on Intelligent Transportation, Big Data and Smart City (ICITBS'15)*, Halong Bay, Vietnam, pp. 81–84.