

# Arabic Multi-Dialect Segmentation: bi-LSTM-CRF vs. SVM

Mohamed Eldesouki<sup>1</sup>, Younes Samih<sup>2</sup>,

Ahmed Abdelali<sup>1</sup>, Mohammed Attia<sup>3</sup>, Hamdy Mubarak<sup>1</sup>, Kareem Darwish<sup>1</sup>, and Laura Kallmeyer<sup>2</sup>

<sup>1</sup>Qatar Computing Research Institute, HBKU, Doha, Qatar

<sup>2</sup>Dept. of Computational Linguistics, University of Düsseldorf, Düsseldorf, Germany

<sup>3</sup>Google Inc., New York City, USA

<sup>1</sup>{mohamohamed, hmubarak, aabdelali, kdarwish}@hbku.edu.qa

<sup>2</sup>{samih, kallmeyer}@phil.hhu.de

<sup>3</sup>attia@google.com

## Abstract

Arabic word segmentation is essential for a variety of NLP applications such as machine translation and information retrieval. Segmentation entails breaking words into their constituent stems, affixes and clitics. In this paper, we compare two approaches for segmenting four major Arabic dialects using only several thousand training examples for each dialect. The two approaches involve posing the problem as a ranking problem, where an SVM ranker picks the best segmentation, and as a sequence labeling problem, where a bi-LSTM RNN coupled with CRF determines where best to segment words. We are able to achieve solid segmentation results for all dialects using rather limited training data. We also show that employing Modern Standard Arabic data for domain adaptation and assuming context independence improve overall results.

## 1 Introduction

Arabic has both complex morphology and orthography, where stems are typically derived from a closed set of roots to which affixes such as coordinating conjunctions, determiners, and pronouns are attached to form words. Segmenting Arabic words into their constituent parts is important for a variety of natural language processing applications. For example, segmentation has been shown to improve the effectiveness of information retrieval (Darwish et al., 2014a) and machine trans-

lation (Habash and Sadat, 2006). Most previous work has mostly focused on segmenting Modern Standard Arabic (MSA) achieving segmentation accuracies of nearly 99% (Abdelali et al., 2016; Pasha et al., 2014). MSA is the lingua franca of the Arab world, and it is typically used in written and formal communications. Dialectal Arabic (DA) segmentation on the other hand has received limited attention, with most of the work focusing on the Egyptian dialect (Habash et al., 2013; Samih et al., 2017). Arabic dialects are typically spoken and are used in informal communications. The advent of the social media and the ubiquity of smart phones has led to a greater need for dialectal processing such as dialect identification (Eldesouki et al., 2016; Khurana et al., 2016), morphological analysis (Habash et al., 2013) and machine translation (Sennrich et al., 2016; Sajjad et al., 2013). Yet, dialectal training corpora for a variety of NLP modules, including segmentation, continue to be limited and often nonexistent.

In this work, we focus on the segmentation of four major Arabic dialects, namely Egyptian, Levantine, Gulf, and Maghrebi. We particularly focus on DA text from Twitter, a popular social media platform, from which we can obtain large amounts of text in different dialects written by ordinary social media users and exhibiting nonstandard orthography. We employ two machine learning approaches for building robust segmentation modules using limited training data (350 tweets containing several thousand words per dialect). In one approach, we pose the segmentation as a ranking problem where all possible segmentations of a word are ranked using a Support Vector Ma-

chine (SVM) based ranker. In the second, we use bidirectional Long Short Term Memory (bi-LSTM) Recurrent Neural Network (RNN) with Conditional Random Fields (CRF) to perform sequence labeling over the characters in words. For both, we adopt the simplifying assumption that word segmentation can be reliably performed independent of context. Though the assumption is not always correct, it has been shown to be fairly robust for more than 99% of word occurrences in Arabic text (Abdelali et al., 2016). Lastly, given the large overlap between MSA and DA, we employ segmented MSA data to further improve dialectal segmentation.

The contribution of this paper are as follows:

- We present robust DA segmenters for four major Arabic dialects. We plan to open-source all of them.
- We provide an exposition of challenges associated with performing in situ DA segmentation including segmentation guidelines and the effect of orthographic standardization.
- We compare two machine learning approaches that can generalize well even when limited training data is available.

## 2 Related Work

Work on dialectal Arabic is fairly new compared to MSA. A number of research projects were devoted to dialect identification (Biadisy et al., 2009; Zbib et al., 2012; Zaidan and Callison-Burch, 2014; Eldesouki et al., 2016). There are five major dialects including Egyptian, Gulf, Iraqi, Levantine and Maghribi. Few resources for these dialects are available such as the CALLHOME Egyptian Arabic Transcripts (LDC97T19), which was made available for research as early as 1997. Newly developed resources include the corpus developed by Bouamor et al. (2014), which contains 2,000 parallel sentences in multiple dialects and MSA as well as English translation.

For the segmentation, Mohamed et al. (2012) built a segmenter based on memory-based learning. The segmenter has been trained on a small corpus of Egyptian Arabic comprising 320 comments containing 20,022 words from [www.masrawy.com](http://www.masrawy.com) that were segmented and annotated by two native speakers. They reported a 91.90% accuracy on the task of segmentation. MADA-ARZ (Habash et al., 2013) is an Egyptian Arabic extension of the Morphological Analysis and Dis-

ambiguation of Arabic (MADA). They trained and evaluated their system on both Penn Arabic Treebank (PATB) (parts 1-3) and the Egyptian Arabic Treebank (parts 1-5) (Maamouri et al., 2014) and they achieved 97.5% accuracy. MADAMIRA<sup>1</sup> (Pasha et al., 2014) is a new version of MADA and includes functionality for analyzing dialectal Egyptian. Monroe et al. (2014) used a single dialect-independent model for segmenting all Arabic dialects including MSA. They argue that their segmenter is better than other segmenters that use sophisticated linguistic analysis. They evaluated their model on three corpora, namely parts 1-3 of Penn Arabic Treebank (PATB), Broadcast News Arabic Treebank (BN), and parts 1-8 of the BOLT Phase 1 Egyptian Arabic Treebank (ARZ) reporting a 95.13% F1 score.

## 3 Dialectal Arabic

### 3.1 DA Challenges

DA shares many MSA challenges, such as having complex templatic derivational morphology and concatenative orthography. Most nouns and verbs are typically derived from a closed set of roots, which are fitted into templates to generate stems. Templates may indicate morphological features such POS tag, gender, and number. Stems may accept prefixes, such as coordinating conjunction and prepositions, or suffixes, such as pronouns, to form words. While dialects mostly comply with the templatic nature of morphology<sup>2</sup>, they diverge from MSA in other aspects such as:

- Lack of standard orthography, particularly for strictly dialectal words such as *عشان* “E\$An” (because), which may also appear as *علشان* “El\$An” or *مشان* “m\$An” (Habash et al., 2012).
- Word borrowing from other languages (Ibrahim, 2006), such as *بلاصتك* “blASTk” (your place) in Maghrebi, or code switching with other languages (Samih et al., 2016).
- Fusing multiple words together by concatenating tokens and dropping letters, such as the word *يقولك* “yqwlk” (he says to you), “yqwl lk” are concatenated and one “l” is dropped.

<sup>1</sup>MADAMIRA release 20160516 2.1

<sup>2</sup>Minor exception exist such the Egyptian template “AtfEI” that occasionally replaces the MSA template “AnfEI” as in *اتكسر* “Atksr” (broke)

- Additional affixes. Dialectal-specific affixes may arise because of: the alteration of pronouns, such as the feminine second person pronoun from كـ “k” to كي “ky” or the plural pronoun تم “tm” to تو “tw”; the introduction of negation prefix-suffix combination ما - ش “mA-\$”, which behaves like the French “ne-pas” negation construct; the placement of present tense markers, such as “b” in Egyptian and Levantine; the use of different future markers such as “H”, “h”, and “g” instead of “s” for MSA; and the shortening of prepositions and fusing them with the words they precede such as the transformation of على “Ely” (on) to ع “E”.
- Letter substitution, where some letters are commonly substituted for others such as “v” which is replaced with “t” in Egyptian (as in كثير “ktyr” (much)) or “q” which is replaced with “j” in Gulf (as in صدج “Sdj” (really)).
- Syntactic differences, such as the use of masculine plural or singular noun forms instead dual and feminine plural, the dropping of some articles and preposition in some syntactic constructs, and the abandonment of some suffixes such as “wn” in favor of “wA” for verbs and “yn” for nouns.

Using raw text from social media introduces additional phenomena such as word elongation, such as أخيرا “>KyyyyrrA” (finally) instead of أخيرا “>KyrA”, and the use of non-Arabic characters such as Urdu characters (Darwish et al., 2012).

## 4 Dataset

We constructed our dataset by obtaining 350 tweets that were authored for each of the following four dialects: Egyptian, Levantine, Gulf, and Maghrebi. For dialectal Egyptian tweets, we obtained the dataset described in (Darwish et al., 2014b), and we used the same methodology to construct the dataset for the remaining dialects. Initially, we obtained 175 million Arabic tweets by querying the Twitter API using the query “lang:ar” during March 2014. Then, we identified tweets whose authors identified their location in countries where the dialects of interest are spoken (ex. Morocco, Algeria, and Tunisia for Maghrebi) using a large location gazetteer (Mubarak and Darwish, 2014). Then we filtered the tweets using a list containing 10 strong dialectal words per dialect, such

as the Maghrebi word كيما “kymA” (like/as in) and the Levantine word هيك “hyk” (like this). Given the filtered tweets, we randomly selected 2,000 unique tweets for each dialect, and we asked a native speaker of each dialect to manually select 350 tweets that are heavily dialectal. Table 4 lists the number of tweets that we obtained for each dialect and the number of words they contain.

Field	Annotation
Orig. word	بيقولك “byqwlk”
Meaning	he is saying to you
In situ Segm.	ب+يقول+ك “b+yqwl+k”
CODA	بيقول لك “byqwl lk”
CODA Segm.	ب+يقول ل+ك “b+yqwl l+k”

Table 1: Egyptian annotation example

Dialect	No of Tweets	No of Tokens
Egyptian	350	6,721
Levantine	350	6,648
Gulf	350	6,844
Maghrebi	350	5,495

Table 2: Dataset size for the different dialects

Segmentation of DA can be applied on the original raw text, or on the cleaned text after correcting spelling mistakes and applying conventional orthography rules, such as CODA (Habash et al., 2012). In this work, we decided to segment the original raw text. Though Egyptian CODA is a reasonably stable standard, CODA for other dialects are either immature or nonexistent. Also, CODA conversion tools are lacking for most dialects<sup>3</sup>. Building such tools requires the establishment of clear guidelines, is laborious, and may require large annotated corpora (Eskander et al., 2013), such as the LDC Egyptian Treebank.

To prepare the ground truth data for a dialect, we enlisted an annotator who is either a native speaker for the dialect or well versed in it and has background in natural language processing. The authors along with another native speaker of the dialect made multiple review rounds on the work of the annotator to ensure consistency and quality. The annotation guidelines were fairly straightforward. Basically, we asked annotators to:

- segment words in a way that would maintain the

<sup>3</sup>except for Egyptian CODA tool that is embedded in MADAMIRA

correct number of part of speech tags

- favor stems when repeated letters are dropped as Table 1
- segment multiple concatenated words with pluses as in the “merged words” example in Table 3.
- attach injected long vowels that trail prepositions or pronouns to the preposition or pronoun respectively (ex. ليكي “lyky” (to you – feminine) → “ly+ky”)
- treat dialectal words that originated as multiple fused words as single tokens (ex. علاش “EIA\$” (why) – originally على أي شيء “EIY >y \$y”)
- do not segment name mentions and hashtags

In what follows, we discuss the advantages and disadvantages of segmenting raw text versus the CODA’fied text with some statistics obtained for the Egyptian tweets for which we have a CODA’fied version as exemplified in Table 1. The main advantage of segmenting raw text is that it doesn’t need any preprocessing tool to generate CODA orthography, and the main advantage of CODA is that it regularizes text making it more uniform and easier to process. We manually compared the CODA version to the raw version of 2,000 words in our Egyptian dataset. We found that in 75.4% of the words, segmentation of original raw words is exactly the same as the their CODA’fied equivalents (ex. ومن “w+mn” (and from) and نعملها “nEml+hA” (we do it)). Further, if we normalize some characters, namely ه ← ي، ة ← ي، ا ← ي، and non-Arabic characters ج ← ج، ح ← ف، گ ← ك، ي ← ي and remove diacritics, the percentage of matching increases to 90.3%. Table 3 showcases the remaining differences between raw and CODA segmentations and how often they appear. The differences are divided into two groups. In the first group (accounting for 6.8% of the cases), the number of word segments remains the same and both the raw and CODA’fied segments would have the same POS tags.

In this group, the “variable spelling” class contains dialectal words that may have different common spellings with one “standard” spelling in CODA. The “dropped letter” and “shortened particles” classes typically involve the omission of letters such as the first person imperfect prefix “>”

Diff.	%	Examples
<b>Same no. of segments and same POS tags</b>		
variable spellings	2.4%	عشان ⇔ عشان “E\$An, EI\$An”
dropped letters	2.3%	ب + حترم ⇔ بها حترم “b+Htrm, b+AHtrm”
merged words	1.4%	يا عم ⇔ يا عم “yA+Em, yA Em”
Shortened particles	0.4%	ع ⇔ علي، ف ⇔ في “E, EIA f, fy”
elongations	0.3%	لييييه ⇔ ليه “lyyyih, lyh”
<b>Different no. of segments or POS tags</b>		
spelling errors	2.2%	انا ⇔ ان ولا ⇔ وإلا “An, Ana wIA, wAlA”
fused letters	0.8%	قالبي ⇔ قال لبي “qAl+y, qAl l+y”

Table 3: Original vs CODA Segmentations

when preceded by the present tense marker ب “b” or the future tense marker ه “h”, and the “A” in negation particle ما “mA” which is often written as م “m” and attached to the following word, and the trailing letters in prepositions. “Merged words” and “word elongations” are common in social media, where users try to keep within limit by dropping the spaces between letters that do not connect or to stress words respectively.

Though some processing such as splitting of words or removing elongations is required to overcome the phenomena in this group, in situ segmentation of raw words would yield identical segments with the same POS tags as their CODA counterparts. Thus, the segmentation of raw words could be sufficient for 97% of words.

In the second group (accounting for 3% of the cases), both may have a different number of segments or POS tags, which would complicate downstream processing such as POS tagging. They involve spelling errors and the fusion of two identical consecutive letters (gemmination). Correcting such errors may require a spell checker. We opted to segment raw input without correction in our reference, and we kept stem, such as verbs and nouns, complete at the expense of other segments such as prepositions as in the example in Table 1.



## 5 Segmentation Approaches

We present here two different systems for word segmentation. The first uses SVM-based ranking ( $SVM^{Rank}$ )<sup>4</sup> to rank different possible segmentations for a word using a variety of features. The second uses bi-LSTM-CRF, which performs sequence-to-sequence mapping to guess word segmentation.

### 5.1 $SVM^{Rank}$ Approach

This approach is inspired by the work done by Abdelali et al. (2016), in which they used SVM based ranking to ascertain the best segmentation for Modern Standard Arabic (MSA), which they show to be fast and accurate. The approach involves generating all possible segmentations of a word and then ranking them.

In training, we generate all possible segmentations of a word based on a closed set of prefixes and suffixes, and the correct segmentation is assigned rank 1 and all other incorrect segmentations are assigned rank 2. Our valid affixes include MSA prefixes and suffixes that we extracted from Farasa (Abdelali et al., 2016) and additional dialectal prefixes and suffixes that we observed during training. Since we are not mapping words into a standard spelling, such as CODA, prefixes and suffixes may have multiple different representations. For example, given the dialectal Egyptian word for “I do not play”, it could be spelled as “m+b+lEb+\$”, “mA+b+lEb+\$”, “mA+b+lEb+\$”, “m+b+lEb+\$y”, “mA+b+lEb+\$y”, etc. In this example, the first prefix could be “m” or “mA” and the suffix could be “\$” or “\$y”.

Here are two example dialectal Egyptian words to demonstrate segmentation:

- Given the input word: عالوش “EAlw\$” (on the face), possible segmentations are: {E+Al+w\$} (correct segmentation), {E+Alw\$}, {E+Al+w+\$}, {E+Alw+\$}, {EAlw+\$}, and {EAlw\$}.
- Given the input word: باديكي “bAdyky” (I give you (feminine)), possible segmentations are: {b+Ady+ky} (correct segmentation), {b+Adyky}, {bAdy+ky}, and {bAdyky}.

We use the following features in training the classifier:

<sup>4</sup>[https://www.cs.cornell.edu/people/tj/svm\\_light/svm\\_rank.html](https://www.cs.cornell.edu/people/tj/svm_light/svm_rank.html)

- Conditional probability that a leading character sequence is a prefix.
- Conditional probability that a trailing character sequence is a suffix.
- probability of the prefix given the suffix.
- probability of the suffix given the prefix.
- unigram probability of the stem (more details about calculating this is showing below).
- unigram probability of the stem with first suffix.
- whether a valid stem template can be obtained from the stem.
- whether the stem that has no trailing suffixes appears in a gazetteer of person and location names (Abdelali et al., 2016).
- whether the stem is a function word, such as على “Ely” (on), من “mn” (from), and مش “m\$” (not).
- whether the stem appears in the AraComLex<sup>5</sup> Arabic lexicon (Attia et al., 2011) or in the Buckwalter lexicon (Buckwalter, 2002). This is sensible considering the large overlap between MSA and DA.
- length difference from the average stem length.

The segmentations with their corresponding features are then passed to the SVM ranker (Joachims, 2006) for training. Our  $SVM^{Rank}$  uses a linear kernel and a trade-off parameter between training error and margin of 100.

Before training the classifier, features needed to be calculated in advance. As training data, we used the aforementioned sets of 350 dialectal tweets for each dialect containing typically several thousand words each. We also use three parts of the Penn Arabic Treebank (ATB); part 1 (version 4.1), 2 (version 3.1), and 3 (version 2), which have a combined size of 628,870 tokens, to lookup MSA segmentations. The intuition behind using such segmented MSA data for lookup is that both MSA and dialects share a fair amount of vocabulary. Thus, using the ATB corpus has the effect of increasing coverage.

We also adopted the simplifying assumption that any given word has only 1 possible correct segmentation regardless of context. Though this assumption is not always true, previous work on MSA has shown that it holds for 99% of the cases (Abdelali et al., 2016). Invoking this assumption has multiple positive implications, namely: we

<sup>5</sup><http://sourceforge.net/projects/aracomlex/>

can use the segmentations that we observed during training directly, which typically cover most common function words, or segmentations that we observed in the ATB, which cover most MSA words that may be prevalent in dialectal text; and we can cache word segmentations leading to significant speedup. Thus, we experimented with three different lookup schemes for every word, namely: 1) we output the rankers guess directly (None); 2) if exists, we use seen segmentations in dialectal training set, and the output of the ranker otherwise (DA); 3) if exists, we use seen segmentation in dialectal training set, else we use segmentation that we observed in the ATB, and lastly the output of the ranker (DA+MSA).

## 5.2 Bi-LSTM-CRF Approach

### 5.2.1 Long Short-term Memory

Recurrent Neural Network (RNN) belongs to a family of neural networks suited for modeling sequential data. Given an input sequence  $x = (x_1, \dots, x_n)$ , an RNN computes the output vector  $y_t$  of each word  $x_t$  by iterating the following equations from  $t = 1$  to  $n$ :

$$\begin{aligned} h_t &= f(W_{xh}x_t + W_{hh}h_{t-1} + b_h) \\ y_t &= W_{hy}h_t + b_y \end{aligned}$$

where  $h_t$  is the hidden states vector,  $W$  denotes weight matrix,  $b$  denotes bias vector and  $f$  is the activation function of the hidden layer. Theoretically, RNNs can learn long distance dependencies, still in practice they fail due vanishing/exploding gradients (Bengio et al., 1994). To solve this problem, Hochreiter and Schmidhuber (1997) introduced the LSTM RNN. The idea consists of augmenting an RNN with memory cells to overcome difficulties with training and efficiently cope with long distance dependencies. The output of the LSTM hidden layer  $h_t$  given input  $x_t$  is computed via the following intermediate calculations: (Graves, 2013):

$$\begin{aligned} i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\ f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\ c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\ o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\ h_t &= o_t \tanh(c_t) \end{aligned}$$

where  $\sigma$  is the logistic sigmoid function, and  $i$ ,  $f$ ,  $o$  and  $c$  are respectively the input gate, forget gate, output gate and cell activation vectors. More

interpretation about this architecture can be found in (Lipton et al., 2015).

### 5.2.2 Bi-directional LSTM

Bi-LSTM networks (Schuster and Paliwal, 1997) are extensions to single LSTM networks. They are capable of learning long-term dependencies and maintain contextual features from past and future. As shown in Figure 1, they comprise two separate hidden layers that feed forward to the same output layer. A bi-LSTM calculates the forward hidden sequence  $\vec{h}$ , the backward hidden sequence  $\overleftarrow{h}$  and the output sequence  $y$  by iterating over the following equations :

$$\begin{aligned} \vec{h}_t &= \sigma(W_{x\vec{h}}x_t + W_{\vec{h}\vec{h}}\vec{h}_{t-1} + b_{\vec{h}}) \\ \overleftarrow{h}_t &= \sigma(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\overleftarrow{h}}\overleftarrow{h}_{t-1} + b_{\overleftarrow{h}}) \\ y_t &= W_{\vec{h}y}\vec{h}_t + W_{\overleftarrow{h}y}\overleftarrow{h}_t + b_y \end{aligned}$$

More interpretations about these formulas are found in Graves et al. (2013a)

### 5.2.3 Conditional Random Fields (CRF)

Over the past few years, bi-LSTMs have achieved many ground-breaking results in many NLP tasks because of their ability to cope with long distance dependencies and exploit contextual features from past and future states. Still, when they are used for some specific sequence classification tasks, (such as segmentation and named entity detection), where there is strict dependence between output labels, they fail to generalize perfectly. During the training phase of the bi-LSTM networks, the resulting probability distributions for different time steps are independent from each other. To overcome the independence assumptions imposed by the bi-LSTM and to exploit these kind of labeling constraints in our Arabic segmentation system, we model label sequence logic jointly using Conditional Random Fields (CRF) (Lafferty et al., 2001).

### 5.2.4 bi-LSTM-CRF for DA Segmentation

In this model we consider Arabic segmentation as a sequence labeling problem at the character level. Each character is labeled with one of five labels  $B, M, E, S, WB$  that designate the segmentation decision boundaries: Beginning, Middle, End of a multi-character segment, Single character segment, and Word Boundary respectively. Figure 1 illustrates our segmentation model and how the model takes the word قلبه “qlbh” (his heart) as its

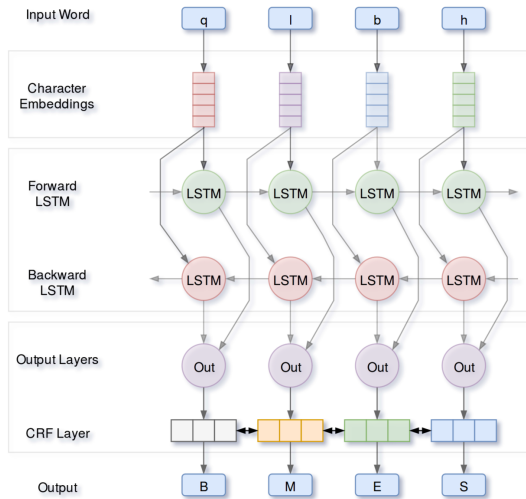


Figure 1: Architecture of our proposed neural network Arabic segmentation model applied to word. قلبه “qlbh” and output “qlb+h”

current input and predicts its correct segmentation. The model is comprised of the following three layers:

- Input layer: containing character embeddings.
- Hidden layer: bi-LSTM maps character representations to hidden sequences.
- Output layer: CRF computes the probability distribution over all labels.

At the input layer, a look-up table is initialized with randomly uniform sampled embeddings<sup>6</sup> mapping each character in the input to d-dimensional vector. At the hidden layer, the output from the character embeddings is used as the input to the bi-LSTM layer to obtain fixed-dimensional representations for each character. At the output layer, a CRF is applied over the hidden representation of the bi-LSTM to obtain the probability distribution over all the labels. Training is performed using stochastic gradient (SGD) descent with momentum 0.9 and batch size 50, optimizing the cross entropy objective function.

### 5.2.5 Optimization

Due to the relatively small size the training and development sets, overfitting poses a considerable challenge for our Dialectal Arabic segmentation system. To make sure that our model learns significant representations, we resort to dropout (Hinton et al., 2012) to mitigate overfitting. The basic

<sup>6</sup>We did not use pre-trained character embeddings, because we conducted side experiments with and without pre-trained embeddings and the results were mixed

idea behind dropout involves randomly omitting a certain percentage of the neurons in each hidden layer for each presentation of the samples during training. This encourages each neuron to depend less on the other neurons to learn the right segmentation decision boundaries. We apply dropout masks to the character embedding layer before inputting to the bi-LSTM and to its output vector. In our experiments, we find that dropout with a fixed rate of 0.5 decreases overfitting and improves the overall performance of our system. We also employ early stopping (Caruana et al., 2000; Graves et al., 2013b) to mitigate overfitting by monitoring the model’s performance on the development set.

## 6 Experiments and Results

As described earlier, we perform several experiments for each dialect. These involve training using dialectal data while using different lookup schemes, namely: no lookup (None); lookup from dialectal training only (DA); and a cascaded lookup from dialectal training and then MSA (DA+MSA). For all our experiments, we use 5 fold cross validation with 70/10/20 train/dev/test splits. We use the Farasa MSA segmenter as a baseline. Table 4 reports on the results for both segmentation approaches and in combination of using different lookup schemes. As the results clearly shows, using an MSA segmenter yields suboptimal results for dialects. Also, when no lookup is used, the bi-LSTM-CRF sequence labeler performs better than the SVM ranker for all dialects. However, using lookup leads to greater improvements for the SVM approach leading to the best results for Levantine, Gulf, and Maghrebi and slightly lower results for Egyptian. Further,  $SVM^{Rank}$  seemed to have benefited more from the DA lookup, while bi-LSTM-CRF benefited more from the MSA lookup. As for Egyptian segmentation, we suspected that it performed better for both approaches than the segmentation for the other dialects, because the percentage of test words that appear in the training set was greater for Egyptian. The percentages for all the dialects are:

Egyptian	Levantine	Gulf	Maghrebi
64.7%	54.7%	56.7%	55.2%

As for the lower results for Maghrebi, we noticed that Maghrebi has many more affixes than MSA and other dialects. These affixes contribute to the data sparsity and complexity of the segmentation task. For example, we enumerated 24 prefixes

Training Set	Look-up	Egyptian	Levantine	Gulf	Maghrebi
SVM <sup>Rank</sup>	None	91.0	87.8	87.7	84.7
	DA	94.5	92.9	92.8	90.5
	DA+MSA	94.6	<b>93.3</b>	<b>93.1</b>	<b>91.2</b>
bi-LSTM-CRF	None	93.8	91.0	89.4	87.1
	DA	94.2	91.8	90.8	88.5
	DA+MSA	<b>95.0</b>	93.0	91.9	90.1
Farasa		85.7	82.6	82.9	82.6

Table 4: SVM<sup>Rank</sup> and bi-LSTM results w/ and w/o lookup

for Maghrebi compared to 8 for MSA, 17 for Levantine and Gulf, and 12 for Egyptian. Similarly, Maghrebi had more suffixes than MSA and other dialects. To ascertain the effect of CODA’fication, we ran an extra experiment where we trained our best SVM<sup>Rank</sup> system using the CODA’fied version of the Egyptian data, and the segmentation accuracy increased from 94.6% to 96.8%. Thus, having stable CODA standards and reliable conversion tools may positively impact dialectal processing. Next, we elaborate on typical errors of both approaches.

**SVM<sup>Rank</sup> Errors:** We examined the errors that the SVM ranker produced for different dialects and the most common types involved:

- erroneous splitting of leading or trailing characters when they were not prefixes or suffixes respectively or not splitting actual prefix and suffixes. For example, “حايكون” “H+ykwn” (will be) was segmented as “Hyk+wn”.
- the use of non-Arabic letters, wrong form of *alef*, or “h” instead of “p”. For example “جاياك” “jAy+l+k” (I am coming to you), where “A” and “k” were replaced with “|” and a Farsi character respectively, was not segmented.
- long words with multiple segments such as “ممتقلقينا” “m+tlq+y+nA+\$” (don’t make us angry) where the ranker chose to segment it as “m+tlq+yn+A\$”.

**bi-LSTM-CRF Errors:** The errors in this system are broadly classified into three categories:

- Ambiguous in token boundary because of character sharing in case of gemination/elision. For example the word “عنا” “En A” (about us) is actually two tokens “En” and “nA”. The two “n” letters are now merged into one. In the gold

data, the disputed letter belongs to the first token while in the system output, it belongs to the second.

- Like the SVM, the system often fails due to unconventional spelling. For example the word “لاخويا” “lAxwyA” (to my brother) is a misspelling of “لأخويا”.
- The majority of the remaining errors are simply mis-tokenization due to the system’s inability to decide whether a substring (which out of context can be a valid token) is an independent token or part of a word, e.g. “مستقبلك” “mstqbl+k” (your future), which is predicted by the system as “m+staqbl+k”, where it correctly recognizes the genitive pronoun in the end, but mistakenly tags the first radical as a separate segment.

## 7 Conclusion

In this paper we presented two approaches involving SVM-based ranking and bi-LSTM-CRF sequence labeling for segmenting Egyptian, Levantine, Gulf, and Maghrebi dialects. Both approaches yield strong comparable results that range between 91% and 95% accuracy for different dialects. To perform the work, we created training corpora containing naturally occurring text from social media for the aforementioned dialects. We plan to release the data and the resulting segmenters to the research community. For future work, we want to perform domain adaptation using large MSA data, such as ATB, to improve segmentation results. Further, we plan to investigate building a joint model capable of segmenting all the dialects with minimal loss in accuracy.

## References

- Ahmed Abdelali, Kareem Darwish, Nadir Durrani, and Hamdy Mubarak. 2016. Farasa: A fast and furious segmenter for arabic. In *Proceedings of the 2016*



- Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*. Association for Computational Linguistics, San Diego, California, pages 11–16.
- Mohammed Attia, Pavel Pecina, Antonio Toral, Lamia Tounsi, and Josef van Genabith. 2011. An open-source finite state morphological transducer for modern standard arabic. In *Proceedings of the 9th International Workshop on Finite State Methods and Natural Language Processing*. Association for Computational Linguistics, pages 125–133.
- Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks* 5(2):157–166.
- Fadi Biadisy, Julia Hirschberg, and Nizar Habash. 2009. Spoken arabic dialect identification using phonotactic modeling. In *Proceedings of the EACL 2009 Workshop on Computational Approaches to Semitic Languages*. Association for Computational Linguistics, Stroudsburg, PA, USA, Semitic '09, pages 53–61.
- Houda Bouamor, Nizar Habash, and Kemal Oflazer. 2014. A multidialectal parallel corpus of arabic. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*. European Language Resources Association (ELRA), Reykjavik, Iceland.
- Tim Buckwalter. 2002. Buckwalter {Arabic} morphological analyzer version 1.0 .
- Rich Caruana, Steve Lawrence, and Lee Giles. 2000. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *NIPS*. pages 402–408.
- Kareem Darwish, Walid Magdy, and Ahmed Mourad. 2012. Language processing for arabic microblog retrieval. In *Proceedings of the 21st ACM international conference on Information and knowledge management*. ACM, pages 2427–2430.
- Kareem Darwish, Walid Magdy, et al. 2014a. Arabic information retrieval. *Foundations and Trends® in Information Retrieval* 7(4):239–342.
- Kareem Darwish, Hassan Sajjad, and Hamdy Mubarak. 2014b. Verifiably effective arabic dialect identification. In *EMNLP*. pages 1465–1468.
- Mohamed Eldesouki, Fahim Dalvi, Hassan Sajjad, and Kareem Darwish. 2016. Qcri@ dsl 2016: Spoken arabic dialect identification using textual features. *VarDial* 3 page 221.
- Ramy Eskander, Nizar Habash, Owen Rambow, and Nadi Tomeh. 2013. Processing spontaneous orthography.
- Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850* .
- Alex Graves, Navdeep Jaitly, and Abdel-rahman Mohamed. 2013a. Hybrid speech recognition with deep bidirectional lstm. In *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*. IEEE, pages 273–278.
- Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013b. Speech recognition with deep recurrent neural networks. In *Acoustics, speech and signal processing (icassp), 2013 IEEE international conference on*. IEEE, pages 6645–6649.
- Nizar Habash, Mona T Diab, and Owen Rambow. 2012. Conventional orthography for dialectal arabic. In *LREC*. pages 711–718.
- Nizar Habash, Ryan Roth, Owen Rambow, Ramy Eskander, and Nadi Tomeh. 2013. Morphological analysis and disambiguation for dialectal arabic. In *Hlt-Naacl*. pages 426–432.
- Nizar Habash and Fatiha Sadat. 2006. Arabic pre-processing schemes for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*. Association for Computational Linguistics, pages 49–52.
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. 2012. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580* .
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Zeinab Ibrahim. 2006. Borrowing in modern standard arabic. *Innovation and Continuity in Language and Communication of Different Language Cultures 9*. Edited by Rudolf Muhr pages 235–260.
- Thorsten Joachims. 2006. Training linear svms in linear time. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, pages 217–226.
- Sameer Khurana, Ahmed Ali, and Steve Renals. 2016. Multi-view dimensionality reduction for dialect identification of arabic broadcast speech. *arXiv preprint arXiv:1609.05650* .
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. ICML*.
- Zachary C Lipton, David C Kale, Charles Elkan, and Randall Wetzell. 2015. A critical review of recurrent neural networks for sequence learning. *CoRR* abs/1506.00019.

- Mohamed Maamouri, Ann Bies, Seth Kulick, Michael Ciul, Nizar Habash, and Ramy Eskander. 2014. Developing an egyptian arabic treebank: Impact of dialectal morphology on annotation and tool development. In *LREC*. pages 2348–2354.
- Emad Mohamed, Behrang Mohit, and Kemal Oflazer. 2012. Annotating and learning morphological segmentation of egyptian colloquial arabic. In *LREC*. pages 873–877.
- Will Monroe, Spence Green, and Christopher D Manning. 2014. Word segmentation of informal arabic with domain adaptation. In *ACL (2)*. pages 206–211.
- Hamdy Mubarak and Kareem Darwish. 2014. Using twitter to collect a multi-dialectal corpus of arabic. In *Proceedings of the EMNLP 2014 Workshop on Arabic Natural Language Processing (ANLP)*. pages 1–7.
- Arfath Pasha, Mohamed Al-Badrashiny, Mona Diab, Ahmed El Kholy, Ramy Eskander, Nizar Habash, Manoj Pooleery, Owen Rambow, and Ryan M Roth. 2014. Madamira: A fast, comprehensive tool for morphological analysis and disambiguation of Arabic. *Proc. LREC*.
- Hassan Sajjad, Kareem Darwish, and Yonatan Belinkov. 2013. Translating dialectal Arabic to English. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Sofia, Bulgaria, ACL '13, pages 1–6.
- Younes Samih, Mohammed Attia, Mohamed Eldesouki, Hamdy Mubarak, Ahmed Abdelali, Laura Kallmeyer, and Kareem Darwish. 2017. A neural architecture for dialectal arabic segmentation. *WANLP 2017 (co-located with EACL 2017)* page 46.
- Younes Samih, Suraj Maharjan, Mohammed Attia, Laura Kallmeyer, and Tamar Solorio. 2016. Multilingual code-switching identification via lstm recurrent neural networks. In *Proceedings of the Second Workshop on Computational Approaches to Code Switching*,. Austin, TX, pages 50–59. <http://www.aclweb.org/anthology/W/W16/W16-58.pdf#page=62>.
- Mike Schuster and Kuldeep K Paliwal. 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 45(11):2673–2681.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, pages 1715–1725. <http://www.aclweb.org/anthology/P16-1162>.
- Omar F Zaidan and Chris Callison-Burch. 2014. Arabic dialect identification. *Computational Linguistics* 40(1):171–202.
- Rabih Zbib, Erika Malchiodi, Jacob Devlin, David Stallard, Spyros Matsoukas, Richard Schwartz, John Makhoul, Omar F. Zaidan, and Chris Callison-Burch. 2012. Machine translation of arabic dialects. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Stroudsburg, PA, USA, NAACL HLT '12, pages 49–59.