# Language Processing and Learning Models
# for Community Question Answering in Arabic

Salvatore Romeo[a], Giovanni Da San Martino[a], Yonatan Belinkov[b],
Alberto Barrón-Cedeño[a], Mohamed Eldesouki[a], Kareem Darwish[a],
Hamdy Mubarak[a], James Glass[b], Alessandro Moschitti[a]

*[a]Qatar Computing Research Institute, HBKU, Doha, Qatar*
*email:* {sromeo, gmartino, albarron, mohamohamed,
kdarwish, hmubarak, amoschitti}@qf.org.qa
*[b]MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA, USA*
*email:* {belinkov, glass}@mit.edu

## Abstract

In this paper we focus on the problem of question ranking in community question answering (cQA) forums in Arabic. We address the task with machine learning algorithms using advanced Arabic text representations. The latter are obtained by applying tree kernels to constituency parse trees combined with textual similarities, including word embeddings. Our two main contributions are: (*i*) an Arabic language processing pipeline based on UIMA —from segmentation to constituency parsing— built on top of Farasa, a state-of-the-art Arabic language processing toolkit; and (*ii*) the application of long short-term memory neural networks to identify the best text fragments in questions to be used in our tree-kernel-based ranker. Our thorough experimentation on a recently released cQA dataset shows that the Arabic linguistic processing provided by Farasa produces strong results and that neural networks combined with tree kernels further boost the performance in terms of both efficiency and accuracy. Our approach also enables an implicit comparison between different processing pipelines as our tests on Farasa and Stanford parsers demonstrate.

*Keywords:* community question answering, constituency parsing in Arabic, tree-kernel-based ranking, long short-term memory neural networks, attention models.

## 1. Introduction

Community-driven question answering (cQA) on the web typically refers to popular forums in which users ask and answer questions on diverse topics. The freedom to post practically any question and answer in virtual anonymity promotes massive participation. The large amount of posts resulting from this environment demands the implementation of automatic models to filter relevant from irrelevant contents. This scenario has received attention from researchers in both the natural language processing and the information retrieval areas. However, for several reasons, languages other than English —including Arabic— have received relatively less attention.

In this research, we focus on the problem of improving the retrieval of questions from an Arabic forum with respect to a new user question. Our task is formally defined as follows. Let $q$ be a new user question and $D$ the set of question–answer pairs, previously posted in a forum. Rank all $\rho \in D$ according to their relevance against $q$. The main purpose of the ranking model is to improve the user's experience by (*i*) performing a live search on the previously-posted questions, potentially fulfilling the user's information need at once and (*ii*) avoiding the posting of similar questions, particularly if they have already been answered. From the natural language processing point of view this can also be the source of a collection of question paraphrases and near-duplicates, which can be further explored for other tasks.

Our model for question ranking uses Support Vector Machines. We use a combination of tree kernels (TKs) applied to syntactic parse trees, and linear kernels applied to features constituted by different textual similarity metrics computed between $q$ and $\rho$. We build the trees with the constituency parser of Farasa —which we introduce in this paper for the first time— and compare it against the well-consolidated Stanford parser [1]. Additionally, we integrated Farasa in a UIMA-based cQA pipeline[1] which provides powerful machine learning features for question similarity assessment and reranking. Furthermore, we design word embeddings to complement the feature vectors.

In contrast to other question-answering (QA) tasks, forum questions tend to be ill-formed multi-sentence short texts with courtesy fragments, context, and elaborations. As TKs are sensitive to long (irrelevant) texts, we focus on the automatic selection of meaningful text fragments to feed TKs. To do so, we design a selection model based on the weights assigned to each word in the texts

---

[1]It should be noted that our UIMA pipeline with Farasa will be made available to the research community.

2

by an attention mechanism in a long short-term memory network (LSTM). Such a model can filter out irrelevant or noisy subtrees from the question syntactic trees, significantly improving both the speed and the accuracy of the TKs-based classifier.

The rest of the paper is organized as follows. Section 2 offers the necessary background on general QA and cQA, both in Arabic and in other languages. In Section 3 we take a brief diversion from QA to describe Farasa, the technology we use for Arabic natural language processing. We turn back to QA in Section 4, where we present our question ranking model. Section 5 describes our neural network model designed to improve our tree representation by selecting the most relevant text fragments. Section 6 discusses our experiments and obtained results. Section 7 concluded with final remarks.

## 2. Background

As models for QA require linguistic resources, work focused on the Arabic language is relatively humble compared to other better-resourced languages, such as English [2]. Obviously, the scarceness of language resources is not the only issue. In Arabic, characteristics such as a rich morphology, the interaction among multiple dialects, and the common lack of diacritics and capitalization in informal language, pose unprecedented challenges for a QA system to succeed [3]. cQA is one specific scenario of QA. Most of the research work carried out for the Arabic language is focused on standard QA: the search for an answer over a collection of free-text documents. Therefore, this section is divided in three parts. Firstly, we overview some of the literature on Arabic QA. Secondly, we describe the three main stages of a cQA system, including a review of the approaches available to tackle each task, mainly for English. Thirdly, we overview the relatively-scarce literature on cQA for Arabic.

### 2.1. Question Answering in Arabic

Here we overview some of the most representative models proposed to address the three components of a QA system in Arabic: question analysis, passage retrieval, and answer extraction.

In *question analysis*, the task consists of generating the best possible representation for a question $q$ in order to retrieve a subset of relevant documents and, eventually, passages. The question pre-processing applied by Rosso et al. [4] consists of stopword removal and named entity recognition. Afterwards, they classify $q$ by means of its intended information need —whether $q$ is asking for a *name*, a *date*, a *quantity*, or a *definition*— in order to look for the required

3

information in the retrieved passages. Other approaches also try to extract the question's focus (i.e., the main noun phrase) as well as named entities [5, 6, 7].

The resulting representation of $q$ is used for retrieving text passages, $p$, that might answer the question. One alternative is retrieving those $p$ that include a certain amount of the words or phrases in $q$. Besides computing a similarity function $sim(q, p)$ [7], the ranking function can be based on the positional distance among the matching terms in the document [8, 9], i.e., the closer the terms in the document, the more likely it may represent a good answer for $q$. A semantic expansion on the basis of resources such as the Arabic WordNet can come into play as well [9].

Once the most promising text passages have been retrieved, it is time to extract specific answers. Most approaches rely on manually-defined patterns, heuristics, rules, and semantic similarities between question focus and candidate answers; for instance, using $n$-grams [6, 10].

By addressing these three generic steps, different kinds of questions can be answered. For instance, Al Chalabi [11] focused on factoid QA by first determining if $q$ is of kind *who*, *what*, *when*, *how*, etc. QASAL (Question-Answering System for Arabic Language) [5] goes beyond factoid QA by exploiting the linguistic annotation system of NooJ [12] to deal with definitional questions as well. Salem et al. [13] focused on *why* and *how* questions by means of the Rhetorical Discourse Structure (RST) formalism.

### 2.2. The Architecture of a Community Question Answering System

The cQA scenario is slightly different: a new question $q$ formulated by the forum user tends to be less factual and more elaborated, often including contextual information, elaborations, multiple questions, and even irrelevant text fragments. The reference collection $D$ is not composed of free-text documents, but of previously-posted forum questions, together with their answers provided by other users (if any). This leads to building a system architecture as the one represented in Figure 1, which is inspired by Potthast et al. [14].

The first step in the cQA architecture is that of heuristic retrieval. Given question $q$ and a relatively-large collection of forum question–answer pairs $\langle \rho, \alpha \rangle \in D$, an inexpensive mechanism is applied to retrieve the most similar (related) questions $\rho$. Standard information retrieval technology (e.g., a search engine based on inverted indexes), can be applied to solve this task. The creators of the corpus [15] we use for our experiments (Section 6) used Solr[2] to deal with
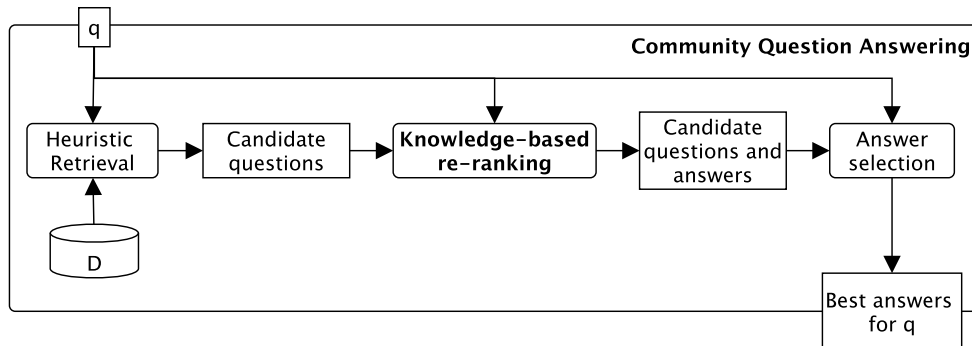
---

[2]https://lucene.apache.org/solr

Figure 1: General architecture of a system for question answering in community-generated forums. *q* stands for the user question; *D* is the collection of previously-posted forum questions along with their answers. The re-ranking stage appears highlighted because it is the problem we address in this research work.

this stage. This step results in the subset of potentially-relevant candidate pairs $D_q \subset D$.

Having $q$ and $D_q$ as input, the *knowledge-based re-ranking* stage is in charge of performing a more refined ordering of the questions. The objective is locating those pairs $\langle \rho, \alpha \rangle \in D$ such that $\rho$ are semantically-equivalent (or at least highly relevant) to $q$. The relatively-small size of $D_q$ allows for the use of more sophisticated —generally more expensive— technology. This is the task we address in this research work, by applying a combination of kernels on both structural and deep learning features (cf. Section 4).

Extensive work has been carried out to design models for this crucial stage of cQA. Although most of them have been devised for English forums, it is worth mentioning some of the approaches. Cao et al. [16] tackled this problem by judging topic similarity, whereas Duan et al. [17] searched for equivalent questions by considering the question's focus as well. Zhou et al. [18] dodged the lexical gap[3] between $q$ and $\rho$ by assessing their similarity on the basis of a (monolingual) phrase-based translation model [19], built on question–answer pairs in a similar fashion to Jeon et al. [20]. Wang et al. [21] computed the similarity between $q$ and $\rho$ on top of syntactic-tree representations: the more substructures the trees have in common, the more similar the questions are. The recent boom in neural network approaches has also impacted question re-ranking. dos Santos et al. [22] applied convolutional neural networks to retrieve semantically-equivalent

---

[3]The classical IR problem of matching the few query terms in relevant documents.

5

questions' subjects. They had to aggregate a bag-of-words neural network when dealing with whole questions; that is, subject and (generally long) body. Support vector machines have shown to be highly competitive in this task. For instance, Franco-Salvador et al. [23] used SVM$^{rank}$ [24] on a manifold of features, including distributed representations and semantic information sources, such as Babel-Net [25] and Framenet [26]. Both Barrón-Cedeño et al. [27] and Filice et al. [28] achieved a good performance using KeLP [29] to combine various kernels with different vectorial and structural features.

Once the most promising questions $\rho$ in the forum are retrieved, potential answers to the new query $q$ are *selected*. The answers $\alpha$ attached to $\rho$ are compared against $q$ in order to estimate their relevance. This is not a trivial problem because the anarchy of Web forums allows users to post irrelevant contents. One of the first approaches to answer selection relied completely on the website's meta-data [30], such as an author's reputation and click counts. Agichtein et al. [31] explored a graph-based model of contributors relationships together with both content- and usage-based features. These approaches depend heavily on the forum's meta-data and social features. Still, as Surdeanu et al. [32] stress, relying on these kinds of data causes the model portability to be difficult; a drawback that disappears when focusing on the content of the questions and answers only. Tran et al. [33] applied machine translation in a similar fashion as Jeon et al. [20] and Zhou et al. [18], together with topic models, embeddings, and similarities. Hou et al. [34] and Nicosia et al. [35] applied supervised models with lexical, syntactic and meta-data features. Some of the most recent proposals aim at classifying whole threads of answers [36, 37] rather than each answer in isolation.

This cQA architecture assumes $q$ is a newly-posted question. A hybrid scenario is that of question deduplication. In this case, $q$ is just another question in the forum, together with its corresponding thread of answers. As a result, the information of both the question and its thread of comments can be used to determine if two posts are asking the same or similar questions. Both Ji et al. [38] and Zhang et al. [39] used LDA topic modeling to learn the latent semantic topics that generate question–answer pairs and used the learned topic distribution to retrieve similar historical questions.

It is worth noting that many of the aforementioned approaches [23, 27, 28, 33, 34, 35] were applied during the two editions of SemEval Task 3 on cQA [40, 15]. In this work we take advantage of the evaluation framework developed for Arabic in the 2016 edition [15] (cf. Section 6.1).

*2.3. Community Question Answering for Arabic*

As the reader can observe, most of the work on cQA has been carried out for other languages than Arabic, including LiveQA [41], which allowed participants to provide answers to *real user questions*, live on the Yahoo! Answers site. To the best of our knowledge, the first effort to come out with a standard framework for the evaluation of cQA models for Arabic is precisely that of [40, 15].

This resource promoted the design of five models for question re-ranking in Arabic. The most successful approach [42] included text similarities at both word and sentence level on the basis of word embeddings. Such similarities were computed both between $q$ and $\rho$, new and retrieved question, respectively, and between $q$ and $\alpha$, with $\alpha$ being the answer linked to the forum question $\rho$ after performing term selection as a pre-processing step. Barrón-Cedeño et al. [27] used tree kernels applied to syntactic trees together with some features in common with [42]. A combination of rule-based, text similarities, and word embeddings has shown to give some benefit in Arabic cQA [43]. Our cQA system reuses ideas and some of the models we developed in [27, 42].

Magooda et al. [44] applied language models enriched with medical terms extracted from the Arabic Wikipedia. Finally, Malhas et al. [45] exploited embeddings in different ways, including the computation of average word vectors and covariance matrices. The performance of these models is included in Table 7, as they represent the state-of-the-art in the testbed we use for our experiments.

## 3. The Farasa Arabic NLP Toolkit

For our Arabic processing, we used our in-house pipeline of Arabic tools called *Farasa*[4] —*insight* or *chivalry* in Arabic. The pipeline includes a segmenter, a POS tagger, a named entity recognizer, a dependency parser, a constituency parser, and a diacritizer. The syntactic parser is a new contribution, introduced in this paper for the first time. Farasa is tuned for the news domain and for Modern Standard Arabic (MSA). Still, Farasa can handle other genres along with classical and dialectal Arabic, but at reduced accuracy. This is possible because of the large overlap between MSA and other varieties of Arabic. Farasa fills an important gap in the span of available tools. It is the only comprehensive suite of Arabic tools that is both open source and whose internal subcomponents are competitive with the state of the art. Here we focus on the relevant components for our current task: segmenter, POS tagger, and constituency parser.
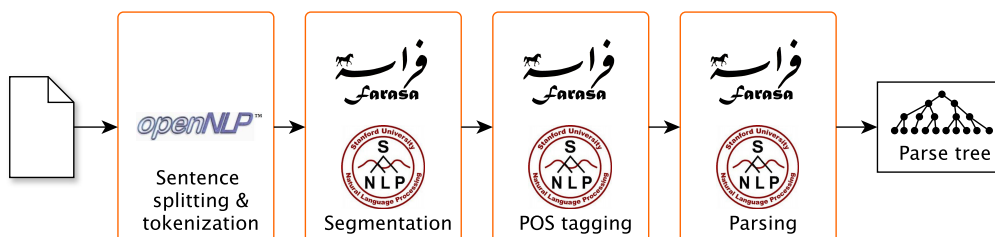
---

[4]Available at `http://farasa.qcri.org`

Figure 2: Our UIMA-based Arabic natural language processing architecture. Each block represents an analysis engine and includes the (alternative) technology it encompasses.

We pose both segmentation and POS tagging as ranking problems, using kernel-based machines. We pose constituency parsing as a sequence labeling problem, where we use a CRF labeler that uses features from the segmenter and POS tagger. Both SVM and CRF have the advantage of being robust and computationally efficient.

### 3.1. UIMA Architecture for Arabic Natural Language Processing

Our Arabic natural language processing pipeline is based on UIMA.[5] UIMA is a framework that allows for the integration of systems to analyze unstructured information (e.g., text documents) whose aim is to extract new knowledge relevant to the particular application context.

UIMA enables to compose applications with self-contained components. Each UIMA component implements an interface defined by the framework and both the input and output structures are described by means of XML descriptor files. The framework is in charge of managing these components, connecting the analysis engines and controlling the data flow. An analysis engine (AE) is a software module that analyzes artifacts (e.g., text) and infers information from them. The analysis engines are built starting from building units called *annotators*. An annotator is a component that analyzes artifacts and produces additional data and/or metadata (e.g., annotation on the analyzed artifact). An AE can contain a single annotator (*primitive AE*) or multiple annotators (*aggregate AE*).

Figure 2 shows the architecture of our pipeline, composed of four AEs. The modularity and flexibility of UIMA allows us for opting for different software modules to perform each of the tasks painlessly. The first AE uses OpenNLP[6]¢¢¢ for sentence splitting, besides performing tokenization. We trained the sentence

---

[5]https://uima.apache.org
[6]https://opennlp.apache.org

8

splitting model on 5*k* sentences from the AQMAR Arabic Wikipedia Supersense corpus [46] and NIST's MT06 corpus.[7] For the rest of the AEs, we can opt for using either Farasa's or Stanford's [1] technology. They are in charge of segmentation into clitics, Part of Speech (POS) tagging, and parsing. In Section 6, we will show the impact of using Farasa or Stanford to process the texts, by comparing different question rankers, each using one of the two parsing systems.

In the following subsections we describe the Farasa segmenter, POS tagger, and parser.

### 3.2. *Farasa Segmenter*

The Farasa segmenter is described in detail in [47, 48]. The segmenter breaks words into their underlying clitics. For example, the word `wktAbhm` (and their book) is segmented into `w+ktAb+hm`. We pose segmentation as a ranking problem, where the ranker attempts to rank possible segmentations of a word. The segmenter uses SVM$^{rank}$ [49] with a linear kernel to determine the best segmentation for each word. We used a linear kernel with a trade-off factor between training errors and margin equal to 100 (parameters tuned on offline experiments carried out over a development set). The ranker uses a dense vector of features which is able to generalize well beyond the cases that are observed during training. Additionally, decoding using SVM$^{Rank}$ is computationally efficient as it involves simple vector multiplication, where speed is highly desirable in processing large amounts of data. We also experimented with using CRF-based sequence labeling [50], and our SVM$^{Rank}$ approach yields better segmentation results with higher speed. Further, we conducted offline experiments to compare our approach to bidirectional Long Short Term Memory (bi-LSTM) over CRF and the results were comparable. It was trained on parts 1 (v. 4.1), 2 (v. 3.1), and 3 (v. 2) of the Penn Arabic Treebank (ATB) [51]. Instead of testing the segmenter on a subset of ATB (which may lead to artificially-high results due to its limited lexical diversity), we tested our segmenter on a corpus of seventy WikiNews articles from 2013 and 2014 [48]. It contains 18, 300 manually-segmented and POS tagged words from articles on seven domains: politics, economics, health, science and technology, sports, arts, and culture.[8]

Table 1 reports on the segmentation accuracy of Farasa and compares it to that of Madamira [52] —a popular state-of-the-art system— on the WikiNews corpus. The performance of the Farasa segmenter is competitive.

---

[7] https://www.nist.gov/programs-projects/machine-translation
[8] The corpus is available at https://github.com/kdarwish/Farasa.

| Task—System | Farasa | Madamira |
|---|---|---|
| Segmentation | 98.9% | 98.8% |
| POS tagging | 94.9% | 95.3% |

Table 1: Accuracy of segmentation and POS tagging for Farasa and Madamira.

### 3.3. Farasa Part-of-Speech Tagger

Our Arabic part-of-speech tagger uses the simplified PATB tag set proposed by [50]. Table 2 includes the tags. The POS tagger attempts to find the optimal tag for each clitic produced by the segmenter, as well as determining the gender (masculine or feminine) and number for nouns and adjectives (singular, dual, or plural). Like the segmenter, the POS tagger uses SVM$^{Rank}$ to find the best tag for each clitic. We decided to adopt SVM$^{Rank}$ for POS tagging for the reasons mentioned earlier for segmentation. Additionally, our SVM$^{Rank}$ outperforms a CRF sequence labeling model [50] and is on par with using a bi-LSTM model [53]. Thus we construct a feature vector for each possible POS tag for each clitic. We supply these vectors to SVM$^{Rank}$ indicating which vector should be ranked the highest given the weights. We then used SVM$^{Rank}$ [49] to learn feature weights. As for the segmenter, we used a linear kernel with a trade-off factor between training errors and margin equal to 100 (parameters tuned on offline experiments carried out over a development set). All possible POS tags for a clitic are scored using the classifier, and the POS with the highest score is picked.

Given a sentence composed of the clitics $c_{-n} \ldots c_0 \ldots c_m$, where $c_0$ is the current clitic and its proposed POS tag, we train the classifier using the following features, computed by maximum-likelihood estimation on our training corpus:

- $p(POS \mid c_0)$ and $p(c_0 \mid POS)$.

- $p(POS \mid c_{-i} \ldots c_{-1})$ and $p(POS \mid c_1 \ldots c_j) \mid i, j \in [1, 4]$.

- $p(POS \mid c_{-i_{POS}} \ldots c_{-1_{POS}})$ and $p(POS \mid c_{1_{POS}} \ldots c_{j_{POS}})$; $i, j \in [1, 4]$. Since we don't know the POS tags of these clitics *a priori*, we estimate the conditional probability as

$$\sum p(POS \mid c_{-i_{possible\_POS}} \ldots c_{-1_{possible\_POS}}) .$$

For example, if the previous clitic could be a NOUN or an ADJ, then $p(POS \mid c_{-1}) = p(POS \mid NOUN) + p(POS \mid ADJ)$.

If the clitic is a stem, we also compute the following features:

10

- $p(POS \mid stem\_template)$. Arabic words are typically derived from a closed set of roots that are placed in so-called stem templates to generate stems. For example, the root `ktb` can be fit in the template `CCAC` to generate the stem `ktAb` (book). Stem templates may conclusively have one POS tag (e.g., `yCCC` is always a verb) or favor one tag over another (e.g., `CCAC` is more likely a NOUN than an ADJ).

- $p(POS \mid prefix)$ and $p(POS \mid suffix)$. Some prefixes and suffixes restrict the possible POS tags for a stem. For example, a stem preceded by DET is either a NOUN or an ADJ.

- $p(POS \mid prefix, prev\_word\_prefix)$, $p(POS \mid prev\_word\_suffix)$ and $p(POS \mid prev\_word\_POS)$. Arabic has agreement rules for noun phrases and idafa constructs (Noun+Noun relation) that cover definiteness, gender, and number. Both these features help capture agreement indicators.

In case we could not compute a feature value during training (e.g., a clitic was never observed with a given POS tag), the feature value is set to $\epsilon = 10^{-10}$. If the clitic is a prefix or a suffix, stem-specific features are assigned the same $\epsilon$ value.

In order to improve efficiency and reduce the choices the classifier needs to pick from, we employ some heuristics that restrict the possible POS tags to be considered by the classifier: (*i*) If the clitic is a number (composed of digits or spelled in words), restrict to "NUM". (*ii*) If all the characters are Latin, restrict to "FOREIGN". (*iii*) If it is a punctuation mark, restrict to "PUNCT". (*iv*) If the clitic is a stem and we can figure out the stem-template, restrict to POS tags that have been seen for that stem-template during training. (*v*) If the clitic is a stem, restrict to POS tags that have been seen during training, given the prefixes and suffixes of the word.

We trained the POS tagger using the same partitions of the ATB that we used for the segmenter (cf. Section 3.2). Table 1 shows the accuracy of our POS tagger on the WikiNews dataset [48] and compares it to Madamira. Madamira edges Farasa by 1.6%. A manual inspection on a random sample of 100 errors showed that 54% of the miss-classifications come from the confusion between adjectives and nouns, whereas 13% are between verbs and nouns. Errors in the preliminary segmentation step cause 21% of the POS mistakes. In such cases, any assigned POS would be incorrect. Table 3 lists the observed error types (covering 95% of errors) including examples.

The POS tagger also assigns gender and number tags to nouns and adjectives. This module is carried over from the Qatara POS tagger [50] and uses the random forest classifier from Weka [54]. The classifier generated 10 trees, with

| POS | Description | POS | Description |
|-----|-------------|-----|-------------|
| ADV | adverb | ADJ | adjective |
| CONJ | conjunction | DET | determiner |
| NOUN | noun | NSUFF | noun suffix |
| NUM | number | PART | particles |
| PREP | preposition | PRON | pronoun |
| PUNC | punctuation | V | verb |
| ABBREV | abbreviation | CASE | alef of tanween fatha |
| FOREIGN | non-Arabic as well as non-MSA words | FUT_PART | future particle "s" prefix and "swf" |

Table 2: Part-of-speech tag set of Farasa.

| Error Type | % | Example |
|------------|---|---------|
| ADJ → NOUN | 29 | "Al<ElAm *Albdyl*" (*alternative* media) "*Albdyl*" recognized as NOUN |
| NOUN → ADJ | 25 | "m$AryE *wykymAnyA*" (*Wikimania* projects) "*wykymAnyA*" recognized as ADJ |
| Segment Error | 21 | "blgp *AlbAyvwn*" instead of "*Al+bAyvwn*" (in *Python* language) |
| V → NOUN | 10 | "hw *Elm* AlErbyp" (he *taught* Arabic) "*Elm*" recognized as NOUN (*science*) |
| Function words | 7 | "mnhA" (from it) recognized as ADJ |
| NOUN → V | 3 | "*k$f* Avry" (archaeological *discovery*) "*k$f*" recognized as V (*discovered*) |

Table 3: POS tagging error types and examples; covering 95% of the errors.

5 attributes for each tree with unlimited depth, and was trained using 8,400 randomly selected unique nouns and adjectives from ATB. The classifier uses the following features: (*i*) stem template; (*ii*) stem template length; (*iii*) POS tag; (*iv*) attached suffix(es); (*v*) whether the word ends with a feminine marker ("At" or "p"); (*vi*) tags that were obtained from a large word list that was extracted from the Modern Arabic Language Dictionary;[9] (*vii*) the 2-gram language-model probability that the word is preceded by masculine or feminine demonstrative articles; and (*viii*) whether the word appears in a gazetteer of proper nouns that have associated gender tags.[10]

For testing, 20-fold cross validation was used. The average accuracy for gender and number classification were 95.6% and 94.9% respectively [50].

## 3.4. Farasa Constituency Parser

The Farasa constituency parser is an in-house re-implementation of the Epic parser [55]; the best-performing Arabic parser in the SPMRL 2013 multilingual constituency parsing shared task [56]. The parser uses a CRF model trained on features derived from the Farasa POS tagger. In compliance with the ATB segmentation, we attached determiners and noun suffixes to the stems. For each clitic, we obtain the information provided by the POS tagger, namely the POS, gender, number, whether the clitic has a determiner, and whether the clitic ends with *ta-marbouta* —the feminine singular noun suffix. Given such information, the parser generates surface features for each clitic $c_0$. Some of these features include the leading and trailing letters in a clitic. The parser uses the leading $n$ letters in the clitic as features ($n \in [1, 5]$). For example, given the clitic `AlktAb` (the book), these features would be {`A,Al,Alk,Alkt,AlktA`}. Similarly, the parser uses the trailing $l$ letters in each clitic as features, ($l \in [1, 5]$). A constraint is placed on the leading and trailing letters: the resulting sequence needs to occur 100+ times in the training data. Furthermore, the parser considers span features, where a span is a bracketed sub-tree (e.g., "(NP (NOUN `AlktAb`))"). The span features include the span's first word, last word, and length; the words before and after the span; split point feature; and span shape feature. To ensure a well-formed nested tree, the parser deduces a minimal probabilistic context-free grammar (PCFG). The parser depends primarily on surface features (i.e. derived only from the clitics in the sentence) to provide context and deep syntactic cues.

---

[9]`http://www.sh.rewayat2.com/gharib/Web/31852/`

[10]We crawled the gazeteer from a list of Palestinian high school graduates including names and genders and Arabic Wikipedia articles (snapshot from September 28, 2012) that have English equivalents and belong to the Wikipedia categories containing the words 'person', 'birth', and 'death' if it has gender information.

|              | POS    | Dev set | Test set |
| ------------ | ------ | ------- | -------- |
| Farasa Parser | golden | **79.70** | 77.01    |
| Farasa Parser | Farasa | 76.94   | 76.34    |
| EPIC Parser   | golden | 78.89   | **78.75** |

Table 4: $F_1$-measure for the Farasa parser compared to the EPIC parser on the SPMRL 2013 shared task dataset. The values are for sentences of all lengths using the *evalb* evaluation script provided by the shared task.

Depending primarily on the surface features gives the parser two advantages. Firstly, it greatly simplifies the structural components of the parser, which would not affect the parser's efficiency since so many deep syntactic cues have surface manifestations. Secondly, it allows for an easy adaptation to new languages.

We used the SPMRL 2013 shared task dataset [57] considering the same training/dev/test partitions for evaluation. In our first experiment, we used the original gold POS tags from the dataset. In our second experiment, we use the segmentation and POS tagging as generated by Farasa. Table 4 compares Farasa (with the two setups) and the Epic parser [55]. Although the Farasa parser is a re-implementation of EPIC, the obtained results differ. Farasa parser when trained with the same dataset as the EPIC parser outperforms it on the dev set, but lags behind on the test with a 1.74 drop in $F_1$ measure. When using the Farasa segmenter and POS tagger to tag words instead of the gold tags we observe a drop of 2.76 and 0.67 for the dev and test sets respectively. The drop can be attributed to tagging errors that are propagated to the parser. However, the drop of 0.67 on the test is an affordable cost for the automation process.

As aforementioned, the Farasa tools are trained on the news genre written in Modern Standard Arabic (MSA), whereas Web forums commonly contain texts written in informal or Dialectal Arabic (DA). Farasa recognizes most of the dialectal words as out of vocabulary (OOV), which affects negatively POS tagging, NER, and syntactic parsing. For a sample of 100 random questions and answers from the Altibbi question-and-answering medical forum,[11] we found that 20% of questions contain at least one dialectal word while answers are written in MSA by professional doctors. In this domain, we found that the majority of the DA words are function words, whereas content words and terms, such as diseases and body parts, are written in MSA. At the semantic level, this is less important compared to the effect at the syntactic level.

---

[11] http://www.altibbi.com; this is the source of the corpus we use in this research.

14

A small degradation in accuracy in Arabic QA systems may occur when using Farasa, designed for MSA, when dealing with DA. Nevertheless, as our results in Section 6 show, this degradation is not important.

## 4. Kernels for Question Re-Ranking

Now we focus on the re-ranking step of cQA, having as input a query question and a set of question-answer pairs, previously retrieved from a Web forum (cf. Section 2.2). Let $Q$ and $\mathcal{A}$ be the set of questions and answers (passages) from the forum, respectively. Let $q$ be a new question. Our task is to model a scoring function, $r : Q \times Q \times \mathcal{A} \rightarrow \mathbb{R}$, which reranks $k$ question–answer pairs, $\langle \rho, \alpha \rangle$, where $\rho \in Q$, $\alpha \in \mathcal{A}$, with respect to their relevance to $q$. Please note that $Q \times \mathcal{A} = D$, which we used in other sections for a more compact reference. We design our scoring function as:

$$r(q, \rho, \alpha) = \vec{w} \cdot \phi(q, \rho, \alpha) \ . \tag{1}$$

We can use implicit representations in kernel-based machines, e.g., SVMs, by expressing $\vec{w}$ as

$$\vec{w} = \sum_{i=1}^{n} \tau_i y_i \phi(q_i, \rho_i, \alpha_i) \ , \tag{2}$$

where $n$ is the number of training examples, $\tau_i$ are weights, $y_i$ are the example labels (*Relevant* and *Irrelevant*), and $\phi(q_i, \rho_i, \alpha_i)$ is the representation of the question pairs. This leads to the following scoring function:

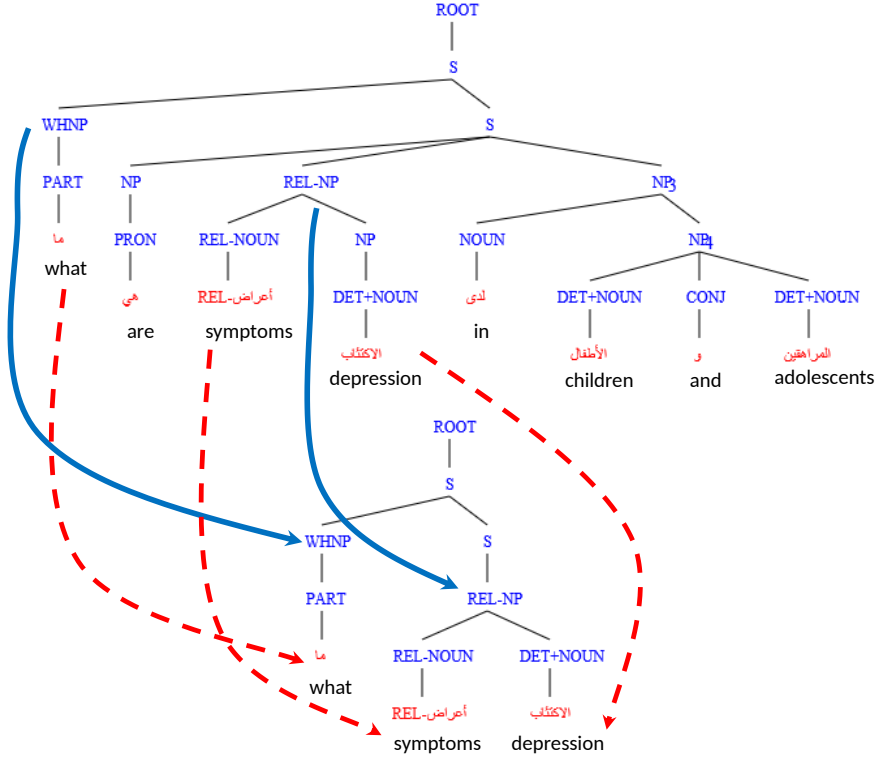$$\begin{aligned} r(q, \rho, \alpha) &= \sum_{i=1}^{n} \tau_i y_i \phi(q, \rho, \alpha) \cdot \phi(q_i, \rho_i, \alpha_i) \\ &= \sum_{i=1}^{n} \tau_i y_i K(\langle q, \rho, \alpha \rangle, \langle q_i, \rho_i, \alpha_i \rangle) \ , \end{aligned} \tag{3}$$

where the kernel $K(\cdot, \cdot)$ intends to capture the similarity between pairs of objects constituted by the query and the retrieved question answer pairs. To any $\phi()$ whose codomain is finite corresponds a kernel function $K(x, x')$, defined on the input space such that $\forall x, x', K(x, x') = \phi(x) \cdot \phi(x')$ [58]. We used three types of representations: parse trees, features derived from word embeddings (word2vec), and text similarity metrics. We combine them as follows:

$$\begin{aligned} K(\langle q, \rho, \alpha \rangle, \langle q_i, \rho_i, \alpha_i \rangle) &= \phi_{tk}(q, \rho) \cdot \phi_{tk}(q_i, \rho_i) && (4) \\ &+ \phi_{w2v}(q, \rho, \alpha) \cdot \phi_{w2v}(q_i, \rho_i, \alpha_i) && (5) \\ &+ \phi_{bow}(q, \rho, \alpha) \cdot \phi_{bow}(q_i, \rho_i, \alpha_i) \ . && (6) \end{aligned}$$

15

$q$: ما هي اعراض الاكتئاب لدى الاطفال والمراهقين
(*What are the symptoms of depression in children and adolescents?*)

$\rho$: ما أعراض الاكتئاب
(*What are depression symptoms?*)

Figure 3: Constituency trees of two questions connected by REL links. The questions correspond to ids 200430 and 47524 in the CQA-MD corpus [15] (cf. Section 6.1).

### 4.1. Tree kernels

We define Eq. (4) as follows

$$\phi_{tk}(q, \rho) \cdot \phi_{tk}(q_i, \rho_i) = TK(t(q, \rho), t(q_i, \rho_i)) + TK(t(\rho, q), t(\rho_i, q_i)) \ , \qquad (7)$$

where $TK$ is a tree-kernel function; e.g., the SubSet Tree (SST) Kernel [59], which measures the similarity between trees. This way, we do not need to extract syntactic feature vectors from the text pairs (i.e., engineering $\phi_{tk}$ is unnecessary). We just need to apply TKs to the pairs of syntactic trees, which provides a score representing the structural similarity. We opt for the state-of-the-art TK model

16

proposed by Severyn and Moschitti [60] and previously used for question ranking in cQA by Barrón-Cedeño et al. [61] and Romeo et al. [62]. As described in Eq. (4), we apply TKs to pairs of questions rather than questions with their answers.

The function $t(x, y)$ in Eq. (7) is a string transformation method that returns the parse tree from the text $x$ —the tree computed with Farasa— further enriching it with the REL tags computed with respect to the syntactic tree of $y$ [60]. The REL tags are added to the terminal nodes of the tree of $x$: a REL tag is added whenever a terminal node of the parse tree of $x$ matches a word in $y$. Typically, REL tags are also propagated to the parent and grandparent nodes (i.e., up to 2 levels). Figure 3 shows the syntactic tree of a query and one of its associated forum questions. The dashed red arrows indicate a matching between words of the two questions, e.g., *Does treatment* or *effect*, whereas the blue arrows are drawn when entire noun phrases or clauses are (partially) matched, i.e., REL-NP or REL-WHNP. The tree nodes are augmented with the REL tag to mark the connection between the constituents of the two syntactic trees.

*4.2. Representation with Embeddings and Similarity Metrics*

Equations (5) and (6) convey a combination of distributional, lexical, and morphosyntactic information from the texts.

To generate the vector $\phi_{w2v}(q, \rho, \alpha)$, we use word vectors obtained with the word2vec tool [63], which is trained (with default settings) on the raw corpus provided with the Arabic cQA task. We compute features that capture similarity between $q$ and $\rho$, and between $q$ and $\alpha$, in the following way. First, we generate a vector representation for every sentence in $q$, $\rho$, and $\alpha$, by averaging the word vectors in the sentence (excluding stopwords). Then, we find the two most similar sentences in $q$ and $\rho$, determined by the cosine similarity between their vector representations, and concatenate their vector representations. We repeat the process for $q$ and $\alpha$ and use their two most similar sentence vectors. Finally, we also find the two most similar word vectors between $q$ and $\rho$ (and between $q$ and $\alpha$), according to the cosine similarity, and add them to the feature representation.

The features in $\phi_{bow}(q, \rho, \alpha)$ from Eq. (6) are obtained using three kinds of text similarity measures applied between $q$ and $\rho$, and between $q$ and $\alpha$: string, lexical, and syntactic. They are included in Table 5.

Our combination of kernels and their corresponding representations is coded in a binary SVM [69].[12] This formulation combines two of the best models presented at SemEval 2016 Task 3 [27, 42, 71] (cf. Section 6.1).

---

[12]Binary SVMs showed comparable results to SVM$^{rank}$ [70].

| Metric | | Details |
|---|---|---|
| String similarity | | |
|     Greedy string tiling | [64] | Considering a minimum matching length of 3. |
|     Longest common subsequence | [65] | Both standard and normalized by the first string. |
|     Longest common substring | [66] | Based on generalized suffix trees. |
| Lexical similarity | | |
|     Jaccard coefficient | [67] | Over stopworded $[1, \ldots, 4]$-grams. |
|     Word containment | [68] | Over stopworded $[1, \ldots, 2]$-grams. |
|     Cosine | | Over stopworded $[1, \ldots, 4]$-grams. |
| | | Over $[1, \ldots, 4]$-grams. |
| | | Over $[1, \ldots, 3]$-grams of part of speech. |
| Syntactic similarity | | |
|     PTK | [59] | Similarity between shallow syntactic trees. |

Table 5: Overview of string, lexical, and syntactic similarity measures.

## 5. Text Selection based on Neural Networks

As shown in Section 2, several neural network approaches have been successfully applied to QA tasks. Unfortunately, question retrieval in cQA is heavily affected by a large amount of noise and a rather different domain, which make it difficult to effectively use out-of-domain embeddings to pre-train neural networks. Figure 4 illustrates some of the difficulties in cQA questions: long greetings and introductions, spelling errors, and incorrect or missing punctuation marks. Correct grammar and usage of punctuation marks is important for sentence splitting and syntactic parsing. This probably prevented the participants to SemEval tasks from achieving satisfactory results with such models [15]. Inspired by [72], in [62] we tried to exploit neural models using their top-level representations for the $(q, \rho)$ pair and fed them into the TK classifier. Nevertheless, this combination proved to be ineffective as well.

Instead of trying to combine the models, we use neural networks to identify the most important pieces of text in both $q$ and $\rho$. We use an LSTM [73, 74], augmented with an attention mechanism. LSTMs have proven to be useful in a number of language understanding tasks. Recently Rocktäschel, et al. [75] adapted an attentional LSTM model [76] to textual entailment, and a similar model has been applied to cQA [77]. We follow the same setup of the latter (Section 5.1). Then, we use the attention weights for our text selection algorithm, which aims at removing subtrees containing useless or noisy information (Section 5.2).

### 5.1. Learning Word Importance with LSTM

The main idea of learning the importance of words for a task is to use the data and labels about the task itself. Given a pair $(q, \rho)$, we learn two serial

18

**Literal Translation:** In the name of God the most beneficent the most merciful our moralist doctor:
I would like to ask you about the <u>bitterness what are its benefits to the body and what are the
benefits and harms of its cholecystectomy</u> please answer as soon as possible and thank you

Figure 4: Example of forum question with long greetings and introductions, spelling errors, and
missing punctuation marks. The most relevant part of the question is underlined.

472 LSTM models: $\text{LSTM}_q$ reads the word vectors of $q$, one by one, and records the
473 corresponding memory cells and hidden states; the final memory cell is used to
474 initialize $\text{LSTM}_\rho$, which reads the word vectors of $\rho$.

Formally, an LSTM computes the hidden representation for input $x_t$ with the
following iterative equations:

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{mi}m_{t-1} + b_i)$$
$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{mf}m_{t-1} + b_f)$$
$$m_t = f_t \odot m_{t-1} + i_t \odot \tanh(W_{xm}x_t + W_{hm}h_{t-1} + b_m)$$
$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{mo}m_t + b_o)$$
$$h_t = o_t \odot \tanh(m_t)$$

475 where $\sigma$ is the sigmoid function, $\odot$ is element-wise multiplication, and $i$, $f$, $o$,
476 and $m$ are input, forget, output, and memory cell activation vectors. The crucial
477 element is the memory cell $m$ that is able to store and reuse long term dependen-
478 cies over the sequence. The $W$ matrices and $b$ bias vectors are learned during
479 training.

480 The final hidden state of $\text{LSTM}_\rho$, $\vec{h}_{\rho,N}$, is used as a feature vector to feed
481 a multi-layer perceptron (MLP) with one hidden layer, followed by a softmax
482 classifier. The objective function is the cross-entropy objective over binary rele-
483 vant/irrelevant target labels.

484 Given the hidden states produced by $\text{LSTM}_q$, we compute a weighted repre-
485 sentation of $q$:

$$\vec{h}_q = \sum_{i=1}^{L} \beta_i \vec{h}_{q,i} \ , \tag{8}$$

486 where $\vec{h}_{q,i}$ are the hidden states corresponding to the words of $q$, and the attention

19

```
 1  Function PruneTree (T, th);
    Input   : a tree T;
              a pruning threshold th;
    Output: a pruned version of T

 2  pruneNode(root(T), th);

 3  Function pruneNode (o, th);
 4  if |children(o)| > 0 then
 5  │   for ch ∈ children(o) do
 6  │   │   pruneNode(ch, th);
 7  │   end
 8  │   if |children(o)| = 0 && !REL_Node(o)) then
 9  │   │   remove (o, T);
10  │   end
11  else
12  │   if o.weight < th && !REL_Node(o)) then
13  │   │   remove (o, T);
14  │   end
15  end
```

**Algorithm 1:** Function *PruneTree* for pruning a tree according to attention weights.

weights $\beta_i$ are computed as:

$$\beta_i = \frac{exp(a(\vec{h}_{q,i}, \vec{h}_{\rho,N}))}{\sum_{j=1}^{L} exp(a(\vec{h}_{q,j}, \vec{h}_{\rho,N}))} \quad . \tag{9}$$

Here $a()$ is parameterized as a MLP with one hidden layer and a *tanh* non-linearity [75]. The input to the MLP is then a concatenation of $\vec{h}_q$ and $\vec{h}_{\rho,N}$.

Intuitively, $\beta_i$ assigns a higher weight to words in $q$ if they are useful for determining the relation to $\rho$. As we will see, these attention weights turn out to be useful for selecting important parts of the questions for the TK models. Note also that the attention here is one-sided —only on $q$. In practice, we train another model, with attention on $\rho$, and use its weights as well.

### 5.2. Parse Tree Pruning based on Neural Networks

Our tree-pruning approach to text selection is illustrated in Algorithm 1. Its main idea is to filter out the leaf nodes of the parse tree corresponding to words

20

associated with weights lower than a user-defined threshold, where the word weights are provided by Eq. (9). The most important step of Algorithm 1 is the recursive function *pruneNode*, which is initially invoked for the root node of the tree. Function *pruneNode* checks whether the node *n* is a leaf (Line 4) and then applies the appropriate strategy: (*i*) for non-leaf nodes, *pruneNode* is invoked for the children of *o*, then *o* is removed if all of its children are removed and (*ii*) a leaf node is removed if its weight is lower than the user-defined threshold, *th*. REL-tagged nodes are never removed, regardless of their weight. Different thresholds determine different percentages of pruned nodes, and we explore various thresholds as part of our experiments.

## 6. Evaluation of Question Re-Ranking Models

In this section, we aim at analyzing the impact of the different representation components in the cQA question re-ranking task. Section 6.1 describes the experimental settings. Section 6.2 illustrates the experimental methodology. Our experiments evaluate four aspects: (*i*) the impact of the NLP processors, (*ii*) the performance of kernels on vectorial features and tree kernels used in isolation, (*iii*) the performance of kernel combinations, and (*iv*) the impact of text selection using tree pruning. We analyze and discuss the results in Section 6.3.

### 6.1. Evaluation Framework

We perform our experiments using the evaluation framework released in the SemEval 2016 Task 3-D [15]. The framework consists of a corpus in Arabic from the medical domain —the CQA-MD corpus— and a set of evaluation metrics. Nakov et al. [15] queried different Web forums to build up a collection of query questions linked to a set of 30 candidate forum questions–answer pairs. The outcome: a total of $45,164$ question–answer forum pairs attached to one of $1,531$ query questions. The relevance of each $\rho \in D$ was manually annotated by means of *crowdsowrcing* considering three labels: *Direct* if $\rho$ contains a direct answer to $q$; *Related* if $\rho$ covers some of the aspects asked by $q$; and *Irrelevant* if $\rho$ and $q$ are unrelated. An ideal ranking should place all direct and relevant $\rho \in D$ on top, followed by the irrelevant pairs. Table 6 shows some statistics of the dataset. The answer associated with each of the 30 forum questions was provided by a professional physician and it is considered correct.

The official evaluation measure is Mean Average Precision (MAP); a standard evaluation metric in information retrieval computed as

$$MAP = \frac{\sum_1^{|Q|} AveP(q)}{|Q|} \ , \qquad (10)$$

| Category | Train | Dev | Test | Total |
|---|---|---|---|---|
| **Questions** | **1,031** | **250** | **250** | **1,531** |
| **QA Pairs** | **30,411** | **7,384** | **7,369** | **45,164** |
| *– Direct* | *917* | *70* | *65* | *1,052* |
| *– Related* | *17,412* | *1,446* | *1,353* | *20,211* |
| *– Irrelevant* | *12,082* | *5,868* | *5,951* | *23,901* |

Table 6: Statistics about the CQA-MD corpus (borrowed from [15]).

where $Q$ is the set of test questions and *AveP* is the average precision value for each query, computed as

$$AveP(q) = \frac{\sum_{k=1}^{|D_q|} (P(k) \times rel(k))}{|\{relevant\ documents\}|} \quad , \tag{11}$$

where $|D_q|$ is the number of retrieved pairs in the ranking, $rel(k)=1$ if $\rho$ at position $k$ is relevant, and $P(k)$ is computed as

$$P(k) = \frac{|\{relevant\ documents\} \cap \{retrieved\ documents\}|_k}{k} \quad ; \tag{12}$$

that is, the size of the intersection between relevant and retrieved documents up to rank $k$ divided by $k$.

### 6.2. *Experiments and Methodology*

Our experiments address the question re-ranking stage in the architecture for community question answering (cf. Section 2). That is, given a query $q$, re-rank a collection of related question–answer pairs in $D_q$. In order to do that, we stick to the same training/development/test partition defined by Nakov et al. [15] for the SemEval 2016 cQA challenge. Regarding the implementation of the models, for the word2vec representations, we trained the embeddings on 26M words of unsupervised data, provided together with the CQA-MD corpus.

We designed four follow-up experiments of increasing complexity:

*Experiment 1: Impact of NLP Processors.* Our first experiment uses only a tree-kernel SVM on parse trees. The difference between our two runs is that we either use Farasa or Stanford's [1] technology to generate the parse-tree representations. This allows for an implicit comparison of these two parsers.

*Experiment 2: Isolated Models.* We perform tests on our three re-ranking models in isolation. Beside the tree-kernel SVM on parse trees from Experiment 1, we experiment with a linear-kernel SVM on word2vec and similarity representations and with the attentional LSTM neural network.

22

| Submission | Dev. | Test |
|---|---|---|
| 1 [42] SLS | 47.31 | 45.83 |
| 2 [27] ConvKN | 42.67 | 45.50 |
| 3 [44] RDI_team | — | 43.80 |
| 4 [45] QU-IR | — | 38.63 |
| 5 [78] UPC_USMBA | — | 29.09 |
| Random Baseline | — | 29.79 |

Table 7: MAP scores of the official submissions to the SemEval 2016 Task 3-D. In addition we report MAP values for the development set of our systems.

*Experiment 3: Kernel Combination.* We combine two SVM kernels on different features: tree kernels on the parse trees and the linear kernel on the word2vec and similarity representations.

*Experiment 4: Tree Pruning.* We explore different thresholds to prune the parse trees on the basis of the LSTM attention weights before learning the scoring function with an SVM. Specifically, we perform experiments combining tree kernels with the linear kernel on word2vec and similarity features.

*6.3. Results and Discussion*

In order to provide a more comprehensive perspective of our experimental results, Table 7 reports the MAP values obtained by the participant systems on the test set of SemEval 2016 Task 3-D. It should be noted that we designed both the two top systems, SLS and ConvKN. The first one was based on a committee of four different systems using different embedding versions as well as methods for filtering the initial word representation, whereas the second applied tree kernels and similarity metrics. In this paper, we only used one system from SLS, corresponding to our linear kernel, which performs relatively more stably with respect to both development and test sets. Although committees are rather effective and typically produce higher accuracy than a single system, they tend to obscure the contribution of the different representations, which are the main target of our study.

It is worth noting that the test set results in Table 7 are obtained by models trained on the training data merged with the development set. Thus, such results are generally higher than those we obtain in this paper on the test set, where we only use the training set for learning all our models. We preferred this approach for our experiments so that we can better compare the results between
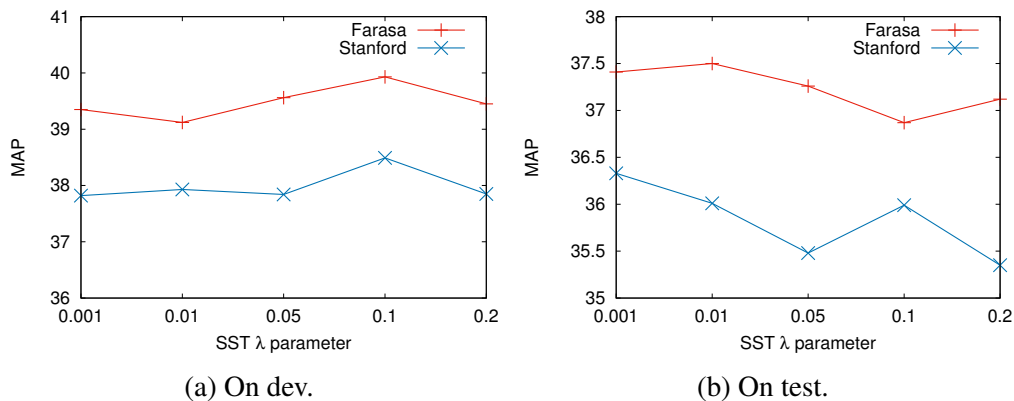
23

Figure 5: MAP as a function of the $\lambda$ parameter of the SST kernel. We compare the performance of our tree-kernel model when the parse-tree representation is built with either Farasa or Stanford.

development and test sets and, at the same time, have a faster training and test processing.

*6.3.1. Experiment 1: Impact of NLP Processors.*

As a way to compare Farasa and Stanford parsers, we ran a set of experiments in which the only difference was the processor used to generate the trees. We used an SVM with $C = 1$ and the normalized SST kernel [79] as TK in Eq. (7) with the following values for the parameter $\lambda = \{0.001, 0.01, 0.05, 0.1, 0.2\}$, which provide different weights to subtrees of different size. Changing $\lambda$, we can emphasize different portions of the parse trees and thus carry out a more systematic comparison between the parsers.

Figure 5 shows the MAP evolution for the two models, with respect to the $\lambda$ parameter of the kernel. The highest MAP values on development (39.93) and test (38.49) sets are obtained when using Farasa. In such cases the increment with respect to Stanford is of 1.44 and 0.88 MAP points, respectively. This is an interesting result as it is in line with our linguistic expert of Arabic who, analyzing some of the trees generated on our data by both parsers, observed a better quality of the Farasa POS-tagger than the one used in the Stanford parser. This different quality also affects chunk definition and their dependencies. It seems that using the entire structure of the parse tree allows TKs to benefit from an overall better quality of Farasa parser to produce better rankings.

24

| Model | Dev. | Test |
|---|---|---|
| Linear-kernel SVM on Word2vec and sims. | 44.94 | 40.73 |
| Tree-kernel SVM on Farasa Parse trees | 42.53 | 40.87 |
| NN (attention on $q$) | 34.85 | 33.40 |
| NN (attention on $\rho$) | 37.47 | 35.09 |

Table 8: MAP performance for our ranking models when applied in isolation on the development and test partitions.

| Model | Dev. | Test |
|---|---|---|
| Tree-kernel (no pruning) + Word2vec and sims. | 46.58 | 41.09 |
| Tree-kernel (pruning ratio 0.74) + Word2vec and sims. | 46.78 | 41.93 |
| Tree-kernel (pruning ratio 0.82) + Word2vec and sims. | 46.01 | 42.20 |

Table 9: MAP performance for our ranking models when applied in combination and after pruning. The latter was applied with two different thresholds, 0.74 and 0.82, which obtained the highest MAP on development and test sets, respectively.

### 6.3.2. Experiment 2: Isolated Models.

Table 8 shows the performance of our ranking models when applied in isolation. The linear- and the tree-kernel models perform on par with each other on the test set, both obtaining competitive results. Still, they lie behind the top 2 systems included in Table 7, at MAP values of $\sim 40.8$ on the test set.

As aforementioned, the neural network does not reach a competitive performance, maybe due to the small amount of data available for training. However, this is not the only contribution the network model can provides as we can use its weights for text selection.

### 6.3.3. Experiment 3: Kernel Combination.

The first row of Table 9 reports the performance of the combination of the tree kernel on parse trees built with Farasa and the linear kernel on word2vec and similarity features. Note that the combination improves over tree kernel and linear kernel in isolation. With respect to our previous systems, i.e., SLS and ConvKN, we got lower values for the test set: as previously pointed out, (*i*) SLS is a combination of four different systems; and (*ii*) in this paper, we only use the training data, whereas we trained SLS and ConvKN on both the training and development sets to obtain the test set results.

### 6.3.4. Experiment 4: Tree Pruning.

While combining feature vectors and tree kernels improves the MAP scores in our experiments, the use of tree kernels has a negative impact on the running time. Thus, we prune parse trees as described in Section 5.2.
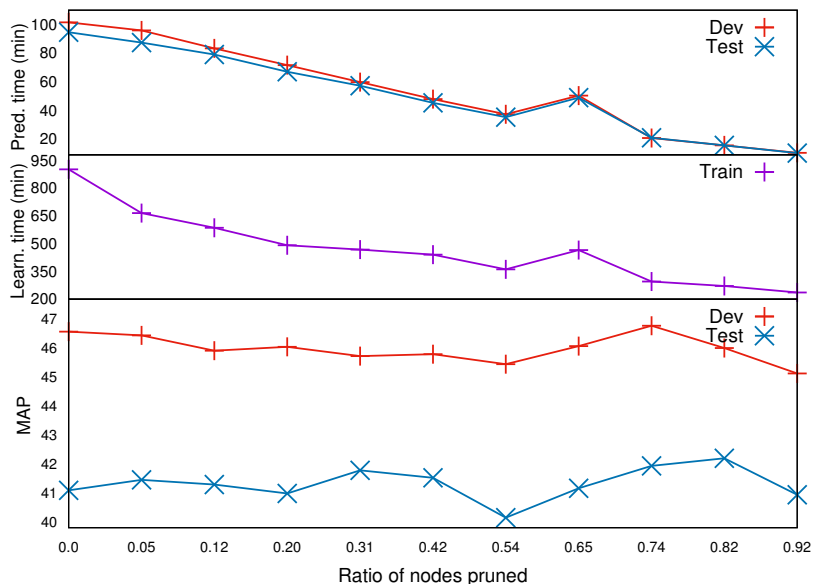
Figure 6: Experiments with pruned trees. From top to bottom the plots show the prediction time, the learning time and MAP as a function of the ratio of pruned nodes.

In this experiment, we evaluate the combination of the linear kernel on word2vec and similarity features with the SST kernel over syntactic trees. Both kernels are not normalized. The top two plots show prediction and learning time (in minutes) as a function of the ratio of pruned nodes. As expected both learning and prediction times decrease roughly linearly with respect to the number of pruned tree nodes.

The plot at the bottom shows the corresponding MAP values, again as a function of the ratio of pruned nodes. Rather than decreasing due to the reduced representation, the MAP scores increase, reaching 46.78 (+0.20 with respect to no pruning) on the development set and 42.20 (+1.11) on the test set. This occurs because our pruning model manages to filter out irrelevant fragments from the trees. For instance, discarding the phrase "in children and adolescents" in Figure 3 would allow a model to better determine that the two questions are practically equivalent.

The threshold maximizing MAP on the development set is the one corresponding to 0.74 pruning ratio (see second line of Table 9). Its MAP score on the test set is 41.93 (+0.84) and the learning and prediction times decrease from 887 to 295 minutes and from 98 to 20 minutes, respectively, with respect to the unpruned data. This means that learning and prediction processes are 3 and 4.9

26

times faster than the kernel combination without pruning.

## 7. Conclusions

Recently, community-driven question answering in websites (cQA) has seen a renewed interest both from natural language processing and information retrieval researchers. Most work in cQA has been carried out for the English language, resulting in a lack of techniques and resources available to deal with other languages, such as Arabic. Motivated by this aspect, in this paper we addressed the problem of cQA in an Arabic forum. In particular, we focused on the task of question re-ranking: given a newly-posted question, retrieve equivalent or similar questions already in the forum. If similar questions have been addressed in the past, the users can quickly obtain an answer to their question.

In order to deal with the necessary processing of the Arabic texts, for the first time, we introduced some components of our in-house pipeline of Arabic NLP tools called Farasa. This includes a segmenter, a POS tagger, a named entity recognizer, a dependency parser, a constituency parser, and a diacritizer. We integrated Farasa into our cQA architecture using the UIMA-based framework. This way, we could extract effective features, such as lexical and syntactic information from Arabic text, and feed them into our machine learning models. Our evaluation on a realistic collection of forum questions in the medical domain allowed us to test Farasa's capabilities when dealing with a real-world application.

In particular, we addressed the task of question re-ranking as a binary classification problem, where each example represents a pair {user-question, forum-question}. We proposed an effective combination of tree kernels built on top of the constituency parse trees provided by Farasa and Arabic word embeddings based on neural networks. This combination allowed for better capturing the semantic relatedness between two short pieces of text, i.e., questions and pairs of questions and answers, and achieved state-of-the-art performance for Arabic question re-ranking.

Additionally, we designed models for selecting meaningful text in order to reduce noise and computational cost. For this purpose, we applied long short-term memory neural networks to identify the best subtrees in the syntactic parsing of questions, which are then used in our tree-kernel-based ranker. We combined the text selection approach with word embeddings based on neural networks, boosting the performance. With thorough experiments we showed that (*i*) syntactic information is very important for the question ranking task, (*ii*) our model combining tree kernels, word embeddings and neural networks for text selection is an effective approach to fully exploit advanced Arabic linguistic processing

27

and (*iii*) our reranker based on tree kernels can be used to implicitly evaluate the performance of different syntactic parsers.

Finally, our UIMA pipeline for Arabic NLP as well as for cQA will be made available to the research community.

## References

[1] S. Green, C. D. Manning, Better Arabic Parsing: Baselines, Evaluations, and Analysis, in: Proceedings of the 23rd International Conference on Computational Linguistics, COLING '10, Association for Computational Linguistics, Stroudsburg, PA, USA, 394–402, URL http://dl.acm.org/citation.cfm?id=1873781.1873826, 2010.

[2] T. Strzalkowski, S. Harabagiu (Eds.), Advances in Open Domain Question Answering, Springer Netherlands, 2008.

[3] A. Ezzeldin, M. Shaheen, A Survey of Arabic Question Answering: Challenges, Tasks, Approaches, Tools, and Future Trends, in: Proceedings of The 13th International Arab Conference on Information Technology (ACIT 2012), 1–8, 2012.

[4] P. Rosso, A. Lyhyaoui, J. Peñarrubia, M. Montes y Gómez, Y. Benajiba, N. Raissouni, Arabic–English Question Answering, in: Proc. of Information Communication Technologies Int. Symposium (ICTIS), Tetuan, Morocco, June, 2005.

[5] M. S. Brini W., Ellouze M., B. H. L., An Arabic Question-Answering System for Factoid Questions, in: IEEE International Conference on Natural Language Procesing and Knowledge Engineering (IEEE NLP-KE'09), Dalian, China, 1–7, 2009.

[6] H. Abdelbaki, M. Shaheen, ARQA High-Performance Arabic Question Answering System, in: Arab Academy for Science and Technology and Maritime Transport, Alexandria, Egypt, 2011.

[7] H. A. Kanaan G., A New Question Answering System for the Arabic Language, American Journal of Applied Sciences 6 (4) (2009) 797–805.

[8] Y. Benajiba, Arabic Question Answering, Master's thesis, Technical University of Valencia, Spain, 2007.

[9] L. Abouenour, On the Improvement of Passage Retrieval in Arabic Question/Answering (Q/A) Systems, in: International Conference on Application of Natural Language to Information Systems, Springer, 336–341, 2011.

[10] O. Trigui, H. Belguith, P. Rosso, DefArabicQA: Arabic Definition Question Answering System, in: Workshop on Language Resources and Human Language Technologies for Semitic Languages, 7th LREC, Valletta, Malta, 40–45, 2010.

[11] H. M. Al Chalabi, Question Processing for Arabic Question Answering System, Ph.D. thesis, The British University in Dubai, 2015.

[12] M. Silberztein, NooJ: a Linguistic Annotation System for Corpus Processing, in: Proceedings of HLT/EMNLP on Interactive Demonstrations, Association for Computational Linguistics, 10–11, 2005.

[13] Z. Salem, J. Sadek, F. Chakkour, N. Haskkour, Automatically Finding Answers to" Why" and" How to" Questions for Arabic Language, in: International Conference on Knowledge-Based and Intelligent Information and Engineering Systems, Springer, 586–593, 2010.

[14] M. Potthast, A. Barrón-Cedeño, B. Stein, P. Rosso, Cross-Language Plagiarism Detection, Language Resources and Evaluation (LRE) 45 (1) (2011) 45–62, ISSN 1574-020X, doi: \bibinfo{doi}{http://dx.doi.org/10.1007/s10579-009-9114-z}.

[15] P. Nakov, L. Màrquez, A. Moschitti, W. Magdy, H. Mubarak, a. Freihat, J. Glass, B. Randeree, SemEval-2016 Task 3: Community Question Answering, in: [80], 525–545, URL http://www.aclweb.org/anthology/S16-1083, 2016.

[16] Y. Cao, H. Duan, C.-Y. Lin, Y. Yu, H.-W. Hon, Recommending Questions Using the Mdl-based Tree Cut Model, in: Proceedings of the 17th International Conference on World Wide Web, WWW '08, ACM, New York, NY, USA, ISBN 978-1-60558-085-2, 81–90, doi:\bibinfo{doi}{10.1145/1367497.1367509}, URL http://doi.acm.org/10.1145/1367497.1367509, 2008.

[17] H. Duan, Y. Cao, C.-Y. Lin, Y. Yu, Searching Questions by Identifying Question Topic and Question Focus, in: [81], 156–164, 2008.

[18] G. Zhou, L. Cai, J. Zhao, K. Liu, Phrase-based translation model for question retrieval in community question answer archives, in: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1, Association for Computational Linguistics, 653–662, 2011.

[19] P. Koehn, F. J. Och, D. Marcu, Statistical Phrase-based Translation, in: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1, NAACL '03, Association for Computational Linguistics, Stroudsburg, PA, 48–54, doi:\bibinfo{doi}{10.3115/1073445.1073462}, URL http://dx.doi.org/10.3115/1073445.1073462, 2003.

[20] J. Jeon, W. B. Croft, J. H. Lee, Finding Similar Questions in Large Question and Answer Archives, in: Proceedings of the 14th ACM International Conference on Information and Knowledge Management, Bremen, Germany, 84–90, 2005.

[21] K. Wang, Z. Ming, T.-S. Chua, A Syntactic Tree Matching Approach to Finding Similar Questions in Community-based QA Services, in: Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval, ACM, 187–194, 2009.

[22] C. dos Santos, L. Barbosa, D. Bogdanova, B. Zadrozny, Learning Hybrid Representations to Retrieve Semantically Equivalent Questions, in: [82], 694–699, 2015.

[23] M. Franco-Salvador, S. Kar, T. Solorio, P. Rosso, UH-PRHLT at SemEval-2016 Task 3: Combining lexical and semantic-based features for community question answering, Proceedings of SemEval 16 (2016) 814–821.

[24] T. Joachims, Training Linear SVMs in Linear Time, in: [83], 217–226, 2006.

[25] R. Navigli, S. P. Ponzetto, BabelNet: The Automatic Construction, Evaluation and Application of a Wide-Coverage Multilingual Semantic Network, Artificial Intelligence 193 (2012) 217–250.

[26] C. F. Baker, H. Sato, The FrameNet Data and Software, in: K. Funakoshi, S. Kübler, J. Otterbacher (Eds.), ACL 2003, 41st Annual Meeting of the Association for Computational Linguistics, Companion Volume to the Proceedings, 7-12 July 2003, Sapporo Convention Center, Sapporo, Japan, The Association for Computer Linguistics, 161–164, URL http://www.aclweb.org/anthology/P03-2030, 2003.

[27] A. Barrón-Cedeño, G. Da San Martino, S. Joty, A. Moschitti, F. Al-Obaidli, S. Romeo, K. Tymoshenko, A. Uva, ConvKN at SemEval-2016 Task 3: Answer and Question Selection for Question Answering on Arabic and English Fora, in: [80], 896–903, URL http://www.aclweb.org/anthology/S16-1138, 2016.

[28] S. Filice, D. Croce, A. Moschitti, R. Basili, KeLP at SemEval-2016 Task 3: Learning Semantic Relations between Questions and Answers, in: [80], 1116–1123, URL http:

29

//www.aclweb.org/anthology/S16-1172, 2016.

[29] S. Filice, G. Castellucci, D. Croce, G. Da San Martino, A. Moschitti, R. Basili, KeLP: a Kernel-based Learning Platform in Java, in: Proceedings of the workshop on Machine Learning Open Source Software: Open Ecosystems, International Conference of Machine Learning, Lille, France, 2015.

[30] J. Jeon, W. B. Croft, J. H. Lee, S. Park, A Framework to Predict the Quality of Answers with Non-Textual Features, in: Proceedings of the 29th ACM SIGIR Conference on Research and Development in Information Retrieval - SIGIR '06, ACM Press, New York, New York, USA, 228, 2006.

[31] E. Agichtein, A. Gionis, C. Castillo, G. Mishne, D. Donato, Finding High-Quality Content in Social Media with an Application to Community-based Question Answering, in: In Proceedings of WSDM, 2008.

[32] M. Surdeanu, M. Ciaramita, H. Zaragoza, Learning to Rank Answers on Large Online QA Collections, in: [81], 719–727, 2008.

[33] Q. H. Tran, V. Tran, T. Vu, M. Nguyen, S. Bao Pham, JAIST: Combining multiple features for Answer Selection in Community Question Answering, in: [84], 215–219, 2015.

[34] Y. Hou, C. Tan, X. Wang, Y. Zhang, J. Xu, Q. Chen, HITSZ-ICRC: Exploiting Classification Approach for Answer Selection in Community Question Answering, in: [84], 196–202, 2015.

[35] M. Nicosia, S. Filice, A. Barrón-Cedeño, I. Saleh, H. Mubarak, W. Gao, P. Nakov, G. Da San Martino, A. Moschitti, K. Darwish, L. Màrquez, S. Joty, W. Magdy, QCRI: Answer Selection for Community Question Answering - Experiments for Arabic and English, in: [84], 203–209, 2015.

[36] S. Joty, A. Barrón-Cedeño, G. Da San Martino, S. Filice, L. Màrquez, A. Moschitti, P. Nakov, Global Thread-level Inference for Comment Classification in Community Question Answering, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Lisbon, Portugal, 573–578, URL http://aclweb.org/anthology/D15-1068, 2015.

[37] G. Zhou, T. He, J. Zhao, P. Hu, Learning Continuous Word Embedding with Metadata for Question Retrieval in Community Question Answering, in: [82], 250–259, 2015.

[38] Z. Ji, F. Xu, B. Wang, B. He, Question-Answer Topic Model for Question Retrieval in Community Question Answering, in: Proceedings of the 21st ACM international conference on Information and Knowledge Management, ACM, 2471–2474, 2012.

[39] K. Zhang, W. Wu, H. Wu, Z. Li, M. Zhou, Question Retrieval with High Quality Answers in Community Question Answering, in: Proceedings of the 23rd ACM international conference on Information and Knowledge Management (CIKM 2014), 371–380, 2014.

[40] P. Nakov, L. Màrquez, W. Magdy, A. Moschitti, J. Glass, B. Randeree, SemEval-2015 Task 3: Answer Selection in Community Question Answering, in: [84], 2015.

[41] E. Agichtein, D. Carmel, D. Harman, D. Pelleg, Y. Pinter, Overview of the TREC 2015 LiveQA Track, in: TREC, 2015.

[42] M. Mohtarami, Y. Belinkov, W.-N. Hsu, Y. Zhang, T. Lei, K. Bar, S. Cyphers, J. Glass, SLS at SemEval-2016 Task 3: Neural-based Approaches for Ranking in Community Question Answering, in: [80], 828–835, URL http://www.aclweb.org/anthology/S16-1128, 2016.

[43] Y. Belinkov, A. Barrón-Cedeño, H. Mubarak, Answer Selection in Arabic Community Question Answering: A Feature-Rich Approach, in: Proceedings of the Second Work-

30

shop on Arabic Natural Language Processing, Association for Computational Linguistics, Beijing, China, 183–190, URL `http://www.aclweb.org/anthology/W15-3223`, 2015.

[44] A. Magooda, A. Gomaa, A. Mahgoub, H. Ahmed, M. Rashwan, H. Raafat, E. Kamal, A. A. Sallab, RDI at SemEval-2016 Task 3: RDI Unsupervised Framework for Text Ranking, in: [80], 822–827, 2016.

[45] R. Malhas, M. Torki, T. Elsayed, QU-IR at SemEval-2016 Task 3: Learning to Rank on Arabic Community Question Answering Forums with Word Embedding, in: [80], 866–871, 2016.

[46] N. Schneider, B. Mohit, K. Oflazer, N. A. Smith, Coarse Lexical Semantic Annotation with Supersenses: An Arabic Case Study, in: Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers - Volume 2, ACL '12, Association for Computational Linguistics, Stroudsburg, PA, USA, 253–258, URL `http://dl.acm.org/citation.cfm?id=2390665.2390726`, 2012.

[47] A. Abdelali, K. Darwish, N. Durrani, H. Mubarak, Farasa: A Fast and Furious Segmenter for Arabic, in: [85], 11–16, 2016.

[48] K. Darwish, H. Mubarak, Farasa: A New Fast and Accurate Arabic Word Segmenter, in: N. C. C. Chair), K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, S. Piperidis (Eds.), Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016), European Language Resources Association (ELRA), Paris, France, ISBN 978-2-9517408-9-1, 1070–1074, 2016.

[49] T. Joachims, Training Linear SVMs in Linear Time, in: [83], 217–226, 2006.

[50] K. Darwish, A. Abdelali, H. Mubarak, Using Stem-Templates to Improve Arabic POS and Gender/Number Tagging, in: [86], 2926–2931, URL `http://www.lrec-conf.org/proceedings/lrec2014/pdf/335_Paper.pdf`, 2014.

[51] M. Maamouri, A. Bies, T. Buckwalter, W. Mekki, The Penn Arabic Treebank: Building a Large-Scale Annotated Arabic Corpus, in: NEMLAR conference on Arabic language resources and tools, vol. 27, 466–467, 2004.

[52] A. Pasha, M. Al-Badrashiny, M. T. Diab, A. El Kholy, R. Eskander, N. Habash, M. Pooleery, O. Rambow, R. Roth, MADAMIRA: A Fast, Comprehensive Tool for Morphological Analysis and Disambiguation of Arabic, in: [86], 1094–1101, 2014.

[53] K. Darwish, H. Mubarak, A. Abdelali, M. Eldesouki, Arabic POS Tagging: Dont Abandon Feature Engineering Just Yet, WANLP 2017 (co-located with EACL 2017) (2017) 130.

[54] L. Breiman, Random Forests, Machine learning 45 (1) (2001) 5–32.

[55] D. L. W. Hall, G. Durrett, D. Klein, Less Grammar, More Features, in: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, 228–237, 2014.

[56] A. Björkelund, Ö. Çetinoğlu, R. Farkas, T. Mueller, W. Seeker, (Re) Ranking Meets Morphosyntax: State-of-the-Art Results from the SPMRL 2013 Shared Task, in: [87] (2013) 135–145.

[57] D. Seddah, R. Tsarfaty, S. Kübler, M. Candito, J. Choi, R. Farkas, J. Foster, I. Goenaga, K. Gojenola, Y. Goldberg, et al., Overview of the SPMRL 2013 Shared Task: Cross-Framework Evaluation of Parsing Morphologically Rich Languages, in: [87], 146–182, 2013.

[58] N. Cristianini, J. Shawe-Taylor, An Introduction to Support Vector Machines and Other Kernel-based Learning Methods, Cambridge University Press, 1 edn., 2000.

[59] A. Moschitti, Efficient Convolution Kernels for Dependency and Constituent Syntactic

          Trees, in: Proceedings of the 17th European Conference on Machine Learning, ECML '06, Springer-Verlag Berlin Heidelberg, Berlin, Germany, 318–329, 2006.

[60] A. Severyn, A. Moschitti, Structural Relationships for Large-scale Learning of Answer Re-Ranking, in: Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '12, Portland, OR, ISBN 978-1-4503-1472-5, 741–750, doi:\bibinfo{doi}{10.1145/2348283.2348383}, URL http://doi.acm.org/10.1145/2348283.2348383, 2012.

[61] A. Barrón-Cedeño, G. Da San Martino, S. Romeo, A. Moschitti, Selecting Sentences versus Selecting Tree Constituents for Automatic Question Ranking, in: [88], 2515–2525, 2016.

[62] S. Romeo, G. Da San Martino, A. Barrón-Cedeño, A. Moschitti, Y. Belinkov, W.-N. Hsu, Y. Zhang, M. Mohtarami, J. Glass, Neural Attention for Learning to Rank Questions in Community Question Answering, in: [88], 1734–1745, 2016.

[63] T. Mikolov, W.-t. Yih, G. Zweig, Linguistic Regularities in Continuous Space Word Representations, in: Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT '13, Atlanta, GA, USA, 746–751, URL http://www.aclweb.org/anthology/N13-1090, 2013.

[64] M. Wise, YAP3: Improved Detection of Similarities in Computer Program and Other Texts, in: Proceedings of the Twenty-seventh SIGCSE Technical Symposium on Computer Science Education, SIGCSE '96, New York, NY, ISBN 0-89791-757-X, 130–134, doi:\bibinfo{doi}{10.1145/236452.236525}, URL http://doi.acm.org/10.1145/236452.236525, 1996.

[65] L. Allison, T. Dix, A Bit-string Longest-common-subsequence Algorithm, Inf. Process. Lett. 23 (6) (1986) 305–310, ISSN 0020-0190, URL http://dl.acm.org/citation.cfm?id=8871.8877.

[66] D. Gusfield, Algorithms on Strings, Trees, and Sequences Computer Science and Computational Biology, Cambridge University Press, 1997.

[67] P. Jaccard, Étude comparative de la distribution florale dans une portion des Alpes et des Jura, Bulletin del la Société Vaudoise des Sciences Naturelles 37 (1901) 547–579.

[68] C. Lyon, J. Malcolm, B. Dickerson, Detecting Short Passages of Similar Text in Large Document Collections, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '01, Pittsburgh, PA, 118–125, 2001.

[69] T. Joachims, Making Large-scale Support Vector Machine Learning Practical, in: B. Schölkopf, C. J. C. Burges, A. J. Smola (Eds.), Advances in Kernel Methods, MIT Press, Cambridge, MA, USA, ISBN 0-262-19416-3, 169–184, URL http://dl.acm.org/citation.cfm?id=299094.299104, 1999.

[70] T. Joachims, Optimizing Search Engines Using Clickthrough Data, in: Proc. KDD, 133–142, 2002.

[71] Y. Belinkov, M. Mohtarami, S. Cyphers, J. Glass, VectorSLU: A Continuous Word Vector Approach to Answer Selection in Community Question Answering Systems, in: [84], URL http://www.aclweb.org/anthology/S15-2038, 2015.

[72] K. Tymoshenko, D. Bonadiman, A. Moschitti, Convolutional Neural Networks vs. Convolution Kernels: Feature Engineering for Answer Sentence Reranking, in: [85], 1268–1278, 2016.

[73] S. Hochreiter, J. Schmidhuber, Long Short-Term Memory, Neural computation 9 (8) (1997) 1735–1780.

32

[74] A. Graves, A.-r. Mohamed, G. Hinton, Speech Recognition with Deep Recurrent Neural Networks, in: Proceedings of ICASSP, 6645–6649, 2013.

[75] T. Rocktäschel, E. Grefenstette, K. M. Hermann, T. Kočiský, P. Blunsom, Reasoning about Entailment with Neural Attention, in: International Conference on Learning Representations, 2016.

[76] D. Bahdanau, K. Cho, Y. Bengio, Neural Machine Translation by Jointly Learning to Align and Translate, arXiv preprint arXiv:1409.0473 .

[77] W. Hsu, Y. Zhang, J. R. Glass, Recurrent Neural Network Encoder with Attention for Community Question Answering, CoRR abs/1603.07044, URL http://arxiv.org/abs/1603.07044.

[78] Y. El Adlouni, I. Lahbari, H. Rodríguez, M. Meknassi, S. O. El Alaoui, UPC-USMBA at SemEval-2016 Task 3: UPC-USMBA participation in SemEval 2016 Task 3, Subtask D: CQA for Arabic, in: [80], 2016.

[79] M. Collins, N. Duffy, Convolution Kernels for Natural Language, in: T. G. Dietterich, S. Becker, Z. Ghahramani (Eds.), NIPS, MIT Press, 625–632, 2001.

[80] Proceedings of the 10th International Workshop on Semantic Evaluation, SemEval '16, Association for Computational Linguistics, San Diego, CA, 2016.

[81] Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics and the Human Language Technology Conference, ACL-HLT '08, Association for Computational Linguistics, Columbus, OH, 2008.

[82] Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, ACL-HLT '15, Association for Computational Linguistics, Beijing, China, 2015.

[83] Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '06, ACM, New York, NY, 2006.

[84] Proceedings of the 9th International Workshop on Semantic Evaluation, SemEval 2015, Association for Computational Linguistics, Denver, CO, 2015.

[85] Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, CA, 2016.

[86] Proceedings of the Ninth International Conference on Language Resources and Evaluation, LREC 2014, European Language Resources Association (ELRA), 2014.

[87] Fourth Workshop on Statistical Parsing of Morphologically Rich Languages, SPMRL '13, Association for Computational Linguistics, Seattle, WA, 2013.

[88] Proceedings of the 26th International Conference on Computational Linguistics, COLING 2016, Osaka, Japan, 2016.

33