

Language Independent Named Entity Recognition via Orthogonal Transformation of Word Vectors

Omar E. Rakha and Hazem M. Abbas

*Dept. Computer and Systems Engineering
Faculty of Engineering
Ain Shams University, Cairo 11571, Egypt*

Abstract

Word embeddings have been a key building block for deep learning NLP in which models relied heavily on word embeddings in many different tasks. In this paper, a model is proposed based on using Bidirectional LSTM/CRF with word embeddings to perform named entity recognition for any language. This is done by training a model on a source language (English) and transforming word embeddings from the target language into word embeddings of the source language by using an orthogonal linear transformation matrix. Evaluation of the model shows that by training a model on an English dataset the model was capable of detecting named entities in an Arabic dataset without neither training or fine tuning the model on an Arabic language dataset.

1. Introduction

Named Entity Recognition (NER) is an important component of information extraction that aims at extracting named entities from textual data. Named entities are words that represent a known person, location, organization or any other entity that can be identified by its name. Extraction of named entities can be used to further improve search engine queries [1] or for question answering [2]. Detection of named entities requires building hand annotated datasets for each target language which consumes significant human labor and time.

NER systems used to rely on hand coded features, part-of-speech (POS) tags, semantic lexicons and huge gazetteers [3]. However this restricted the model's performance and required huge work and domain knowledge to find suitable features.

Current state-of-the-art systems in NER are very accurate with performance exceeding 90% [4] [5] by using a neural network composed of a Bidirectional LSTM layer [6] and a Conditional Random Field (CRF) classifier [7]. However, these models are only trained for a specific language (English for example) and can not be used on any other language.

Due to language dependence, extending the knowledge obtained by these models to new languages is impossible. Thus trying to support new languages requires building new models that only work on the language they are trained on. The effect of this is that the number of languages that can be supported in an application is limited by the amount of labeled data that can be obtained in every language and the performance of each model.

In this work, a novel method is proposed for NER by training a model on an English corpus and evaluating the resulting model on an Arabic corpus by transforming

25 the Word embedding space from Arabic to English. The method can be extended
 26 to any language as long as it is possible to create the transformation matrix for its
 27 word embeddings. The proposed model does not explicitly rely on language specific
 28 features since the employed word embedding model implicitly captures language
 29 specific features and benefits from morphological features of words. This brings
 30 about a model that is both simple and efficient.

31 The paper organization is as follows. The problem is defined in details in Sec-
 32 tion 2. Previous NER models are reviewed in Section 3. The proposed language
 33 independent NER model is presented in Section 4. Experimental results produced
 34 when the model is applied to a new target language are analyzed in Section 5 and
 35 Section 6 concludes the paper.

36 2. Problem Statement

37 The language independent NER problem can be composed of two sub problems
 38 the first being the detection of named entities and the second being language in-
 39 dependence. The sub problem of detecting named entities is a sequence labeling
 40 problem in which the goal is to find out words that are named entities (e.g.,: Per-
 41 sons, Organizations, Locations, ...).

42 A named entity is a word or a phrase that stands consistently for some referent.
 43 This includes names of people, places, organizations. This can also include temporal
 44 expressions (e.g., January fifth 2010, August, 20 November). Numerical expressions
 45 can also be considered named entities (e.g., \$50, 25.5%).

46 For example in Figure 1 each word is labeled in a text sequence as a named
 entity (In this case a person and a location). Outside words are labeled with "O"

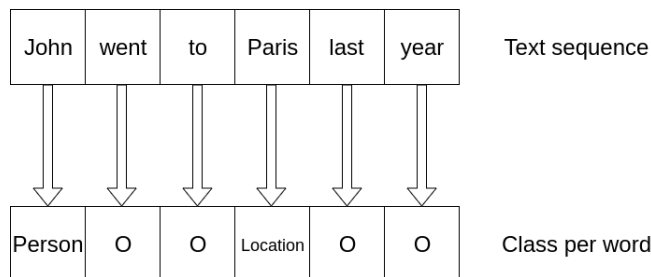


Figure 1: An NER Example

47 following the IOB (Inside-Outside-Beginning) format presented in [8]. The IOB
 48 format prefixes the class of the beginning word with a "B-" indicating it's the first
 49 word in this named entity. while the following words that belong to the same entity
 50 will be prefixed with "I-" which stands for inside. This way we can label named
 51 entities that are longer than 1 word. Refer to Figure 2 for illustration.

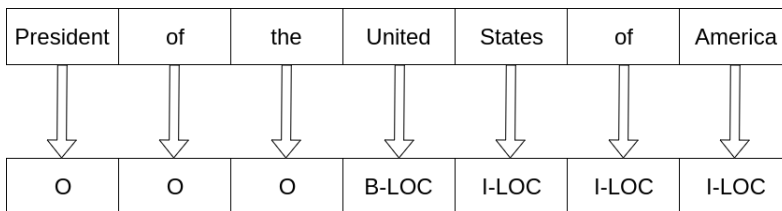


Figure 2: The IOB format

52

53 The second sub problem addressed here is language independence. The goal of
54 building language independent models is to have solutions that can ultimately solve
55 problems regardless of the language it receives. Ideally these models will learn the
56 underlying structure of human language and thus be able to tackle the problems
57 addressed at it without being concerned with the language it uses.

58 For example there exist many data sets for NER in English. However, Arabic
59 data is very scarce. A model that is trained using English data and can be extended
to solve Arabic input is said to be language independent (Figure 3).

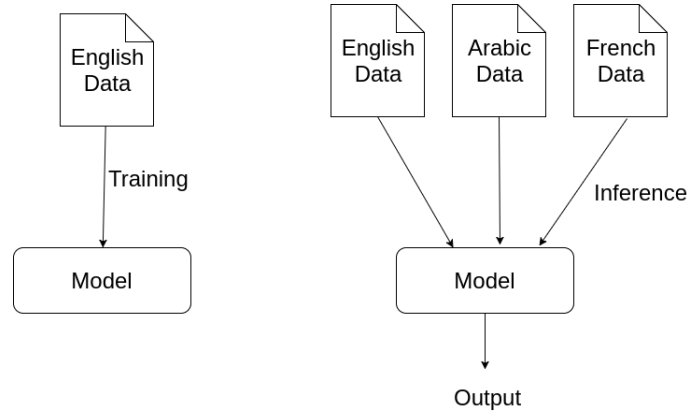


Figure 3: A language independent model

60

61 3. Related Work

62 Currently NER multilingual systems approach the problem by relying on data
63 sets that are designed for building multilingual models like the UN parallel corpora
64 [9]. In following subsections we will explore these approaches and their limitations
65 before proposing the new approach.

66 3.1. Cross Lingual Resources

67 The cross lingual resources approach is centered around finding common features
68 in languages or features that are language dependent. The goal of this approach is
69 to unify the features of more than one language in order to build a classifier that is
70 capable of finding the named entities based only on these features.

71 The system [10] incorporates Wikipedia and other knowledge bases like DBpedia
72 [11]. It uses features like cross-lingual capitalization, transliteration and DBpedia
73 based labeling.

74

75 **Cross-lingual Capitalization:** Since Arabic has no capitalization DBpedia
76 was used to find capitalization weights for a word and be used as a feature.

77 An example: the phrase "Pacific Ocean" was capitalized 36.7% of the time.
78 Thus, the arabic word for "Ocean" was assigned the feature "B-CAPS-0.4" and the
79 arabic word for "Pacific" was assigned the feature "I-CAPS-0.4". The prefix of the
80 feature determines if the word was a "Beginning" or "Inside" word. The body of the
81 feature "CAPS" symbols capitalization and the suffix of the feature is the frequency
82 of being "CAPS".

83 **Transliteration:** The intuition behind transliteration is that named entities are
84 usually transliterated rather than translated. For example the word "Hassan" is an
85 arabic name and also means "good". They use a transliteration weight as a feature
86 similar to the capitalization weight.

87 **DBpedia based labeling:** DBpedia is a large collaboratively-built knowledge
88 base in which structured information is extracted from Wikipedia. For exam-
89 ple, NASA is assigned the following types: Agent, Organization, and Government
90 Agency. These types were used as a feature for each word (if it exists in DBpedia).

91 The problem of this approach is the heavy reliance on the nature of a language,
92 i.e., it has capitalization, transliteration is valid, and it is well represented in DBpe-
93 dia.

94 *3.2. Multi Task Cross Lingual Training*

95 Due to the limitations of the previous approach in having to find features that
96 are viable in more than one language, like Capitalization, a better approach was
97 needed to overcome this limitation. A suitable approach was to rely on embedding
98 to let the model find these features by itself.

99 In multi task cross lingual training [12], a model is trained on multiple tasks such
100 as POS, chunking and NER for more than one language at the same time by using
101 a character level embedding GRU [13] model and a word level embedding model.
102 This model is trained with a shared embedding layer to capture the morphological
103 similarity between languages and learn word vectors that are based on the two
104 languages. This embedding space will capture the similarities between words from
105 the two languages. The final layer is a Conditional Random Field (CRF) classifier
106 [7] which has specific weights for each task.

107 The model in (Figure 4) is structured as pluggable components. The bottom
108 component is an embedding layer. This layer is shared among different tasks and
109 different languages, which means that the model will use the same embedding layer
110 if it was trained for English NER, English POS, Arabic Chunking or Spanish POS.
111 The preceding layer is a task dependent CRF [7] classifier that works on top of the
112 shared embedding layer to produce task specific results (for example: named entities
113 in case of NER).

114 However this model will require a huge amount of annotated data in every needed
115 language. In addition, the model will require retraining whenever more languages
116 need to be added. Also the parameters of the model will increase significantly with
117 the number of languages and hence increase the complexity of the model. The shared
118 embedding layer will not be able to capture any similarity between languages that
119 are intrinsically different like English and Arabic.

120 These approaches show promising results on the multilingual NLP tasks, however
121 the growing complexity of these models with the number of languages makes it near
122 impossible to have models that can work on a huge number of languages. Build-
123 ing a model to work with a few languages would consume a considerable amount
124 of resources(time, annotation effort, ...). Therefore an approach that would work
125 independently of the number of languages is needed.

126 **4. Proposed Approach**

127 In this section, an approach to build a model that is both efficient and can
128 scale reliably with the number of languages presented. The idea of this approach

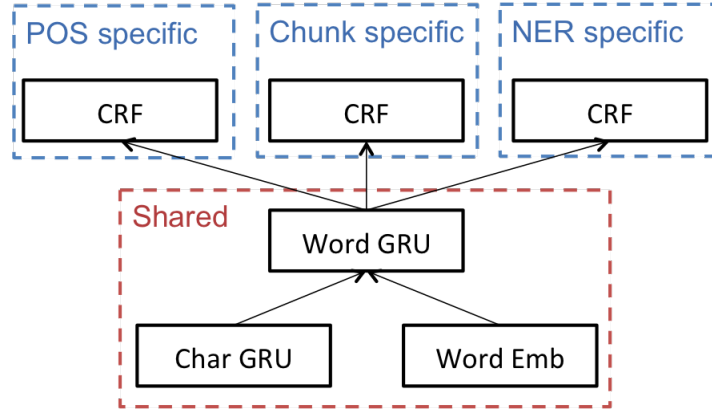


Figure 4: The Multi Task Cross Lingual Training model

129 is to build a neural NER model using the standard components for the problem:
 130 Word embedding for embedding words into vectors, Bidirectional LSTM for feature
 131 extraction and a Conditional Random Field sequence classifier.

132 However to introduce language independence to the model a new step; Word
 133 embedding transformation, is added to the pipeline, which will be able to transform
 134 the vectors of the words of a certain target language (Arabic in our case) into the
 135 space of the vectors of the source language (English in our case) which we trained our
 136 model on. This way the model will be independent of the language of the input and
 137 will only depend on the quality of the Word Embedding Transformation technique
 138 that is applied. Thus the model will be composed of four components:

- 139 1. Word Embedding
- 140 2. Word Embedding Transformation
- 141 3. Bidirectional LSTM Layer
- 142 4. Conditional Random Field

143 Figure 5 shows the architecture of the proposed neural architecture. First, each
 144 word in a sentence is embedded into its corresponding word vector using the word
 145 embedding model. If the word vectors are in a language different than the source
 146 language (English in this case) the word embedding transformation is invoked to
 147 transform the vector from the space of the target language to the space of the
 148 source language. The resulting word vector sequence is then used as the input to a
 149 bidirectional LSTM layer which encodes the input sequence. Finally a CRF layer is
 150 used to tag the sequence.

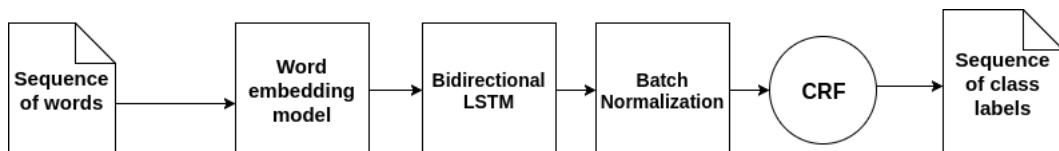


Figure 5: Architecture of the proposed model

151 4.1. Word Embedding

152 Because language is a sequence of words, while neural networks only work with
 153 numbers we have to transform our words into a numerical format. Vector space

154 models are the numerical representation of documents in which a vector is used as a
 155 representation for this document. An example of this is Bag of Words(BOW) [14].

156 In BOW models a vector is assigned for each document representing word counts
 157 for each word in the vocabulary in this document. Consider our vocabulary to only
 158 have the words: ["the", "boy", "girl", "rides", "bus", "and", "from"]. A representa-
 159 tion of the string "the boy rides the bus and the girl rides the bus" would be [4, 1,
 160 1, 2, 2, 1, 0] which is the count of each word in the vocabulary found in this string.

161 Vectors for words can also be computed using BOW by having a one hot repre-
 162 sentation over the vocabulary. For example the word "boy" would receive the vector
 163 [0, 1, 0, 0, 0, 0, 0]. this way we can have vectors representing every word in our
 164 vocabulary.

165 The limitation of this approach is that vectors are very sparse since for a word
 166 the vector will be mostly zeros except for a single value that corresponds to the
 167 word. These vectors also lack any information on the semantic meanings of these
 168 words since the position of the word in the vocabulary gives no information about
 169 similar words to it.

170 Due to these limitations distributed representation of words outperforms tradi-
 171 tional BOW methods because they are dense vectors that can learn semantic and
 172 syntactic properties of these words.

173 A famous example for this taken from [15] is how their model learned a rela-
 174 tionship between Country-Capital in the PCA projection of these vectors (Figure
 175 6).

176 From the 2 dimensional PCA projection we can see that the model learned sim-
 177 ilar representation for countries (China, Russia, Japan, ...) and similar representa-
 178 tion for capitals (Beijing, Moscow, Tokyo, ...) and the direction of the vector be-
 179 tween a country-capital is similar. For example the relationship $model("China") -$
 180 $model("Beijing") \approx model("Russia") - model("Moscow")$ can be seen from the
 181 projection.

182 These models are trained on huge corpora of data -Wikipedia for example- and
 183 the model learns from the context information of words the best representation that
 184 would capture semantic and syntactic information and embed it into the highly
 185 dense vector.

186 Our model employs the Word embedding model introduced in [16] which is a Skip
 187 Gram model [15] with subword information. Given a matrix \mathbf{W} , an entry \mathbf{w}^i is the
 188 i th vector in the matrix which is the vector representation of the word i . The Skip
 189 Gram's objective is to find word representations that are useful for predicting the
 190 surrounding words in a sentence or a document. The model is trained to maximize
 191 the average log probability over a sequence of words $\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3, \dots, \mathbf{w}_T$ where c is
 192 the size of the training context (training window) and \mathbf{w}_t is the center word.:

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(\mathbf{w}_{t+j} | \mathbf{w}_t) \quad (1)$$

193 The basic skip gram model defines the probability of word $t + j$ given word t as
 194 $p(\mathbf{w}_{t+j} | \mathbf{w}_t)$ using the softmax function:

$$\log(p(\mathbf{w}_O | \mathbf{w}_I)) = \log\left(\frac{\exp(v_{\mathbf{w}_O}^\top v_{\mathbf{w}_I})}{\sum_{w=1}^W \exp(v_w^\top v_{\mathbf{w}_I})}\right) \quad (2)$$

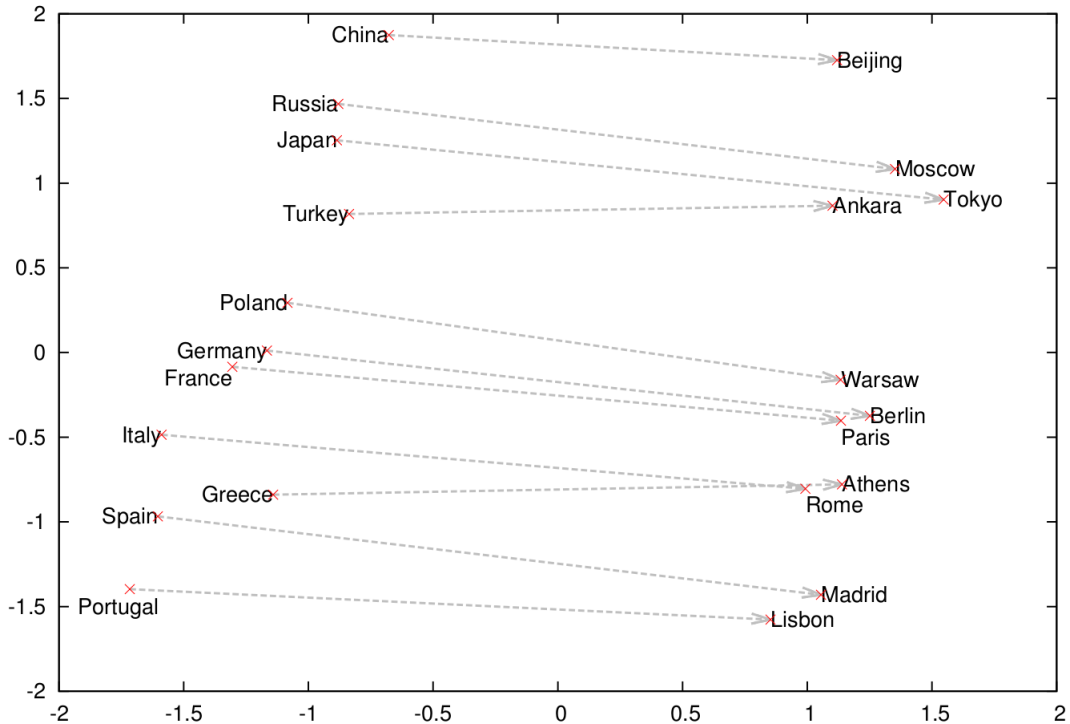


Figure 6: Country-Capital Relationship

195 where v_w and v'_w are the input and output vector representations of \mathbf{w} and W is the
 196 number of words in the vocabulary. This can be simplified into:

$$v'_{\mathbf{w}_O} v_{\mathbf{w}_I} - \log\left(\sum_{w=1}^W \exp(v'_w v_{w_I})\right) \quad (3)$$

197 which consists of two terms, the dot product $v'_{\mathbf{w}_O} v_{\mathbf{w}_I}$ which represents the similar-
 198 ity between the two vectors and the other term $\log(\sum_{w=1}^W \exp(v'_w v_{w_I}))$ being the
 199 normalization term. Computing the normalization term is very expensive since it
 200 requires summing over all words in the vocabulary. This makes it impractical to
 201 compute the gradient of the softmax term $\nabla p(\mathbf{w}_O | \mathbf{w}_I)$ which is proportional to the
 202 size of the vocabulary W .

203 Negative sampling (NEG) is used as the objective for training the Skip Gram
 204 model instead of the softmax to overcome the problem of computing the normaliza-
 205 tion factor. NEG is defined as:

$$\log \sigma(v'_{\mathbf{w}_O} v_{\mathbf{w}_I}) + \sum_{i=1}^k \mathbb{E}_{\mathbf{w}_i \sim P_n} [\log \sigma(-v'_{\mathbf{w}_i} v_{\mathbf{w}_I})] \quad (4)$$

206 Negative sampling works by randomly selecting a set of k "negative" words drawn
 207 from a uni-gram distribution (P_n), These words are the words the model should not
 208 be predicting and instead of updating the weights for every word in the vocabulary
 209 that is not the target word, the optimizer will only update weights for this negative
 210 sample avoiding a huge computation that would make model training in-feasible.

211 In other words the vectors are trained by teaching the model to differentiate
 212 between the target word and noise.

213 So far the model is only capable of learning a distinct vector representation per
 214 word. However, this ignores the morphological features of words (prefixes, suffixes,
 215 n-grams, ...). Therefore, the scoring function needs to be adjusted to take this
 216 information into account.

217 A proper modification for the model would be to learn vector representations
 218 for words and for character n-grams instead of learning vectors for words only. The
 219 implication of this is that each word will be decomposed into a set of n-grams and
 220 the word itself, each element of this set has a distinct vector. this allows for n-gram
 221 vector sharing between words thus producing high quality vectors that take into
 222 account the structure of the word itself.

223 For $\varrho_w \subset 1, \dots, G$ with G being the set of n-grams appearing in \mathbf{W} and ϱ_w being
 224 the set of n-grams appearing in word w , a vector representation is associated with
 225 each n-gram, g , a word is represented by the sum of the vectors of its n-grams and
 226 a unique vector for the word itself, the scoring function is therefore:

$$s(w, c) = \sum_{g \in \varrho_w} z_g^\top v_c \quad (5)$$

227 Where z_g is a vector representation of each n-gram g .

228 For example the word "Good" can be represented by the set of character n-grams
 229 of length 3,4 "Goo", "ood", "Good" where every n-gram of which will have its own
 230 vector representation.

231 4.2. Word Embedding Transformation

232 The idea of this model is to be trained on a corpus from a source language
 233 (English) Figure 7a -since there are available annotated data- and test the model
 234 on a target language with less data (Arabic) Figure 7b, this wouldn't be possible if
 235 we use the word embedding of the two languages separately. This is due to them
 236 having different distributions and the model will not understand the target word
 237 embedding. However, we follow an approach [17] to transform the word embedding
 238 of our target language to our source language and then use it on our model Figure
 239 7c. This approach is tested by training the model on an English dataset [18] and
 240 then evaluate the model on an Arabic dataset [3].

241 We can see in Figure 7a two English words projected on a 2 dimensional space
 242 and in Figure 7b the same two words in Arabic projected on their embedding space.
 243 We can clearly see that the two embedding spaces are different and that the two
 244 words do not overlap. The goal of Word Embedding Transformation is to find a
 245 feasible transformation that can transform the vectors of one of the two spaces into
 246 the other space in order to have words with the same meaning ideally overlapping
 247 Figure 7c.

248 The transformation can be performed by using a translation matrix that is ca-
 249 pable of approximating target word embedding into source word embedding. Using
 250 a translation matrix is a powerful approach that will map one matrix to another by
 251 matrix multiplication.

252 Suppose we have a set of word embeddings that correspond to our target and
 253 source languages X, Y where $X \in R^{n*d}$ and $Y \in R^{n*d}$ such that n is the number of
 254 vectors and d is the dimension of each vector, our goal is to find a matrix W that
 255 approximates the function:

$$\sum_{i=1}^n ||Wx_i - y_i||^2 \quad (6)$$

256 Where x_i is our target language vector and y_i is the source language vector for
 257 the i th word. We can solve this equation using Stochastic Gradient Descent to find
 258 the optimal matrix W that minimizes the error function. Now by using this matrix
 259 W we can translate any vector from our target language X into our source language
 260 Y by multiplication with Wx_i .

261 However, this approach produces sub optimal results due to falling to local min-
 262 ima and being susceptible to learning rates, as well as being time consuming as you
 263 have to iterate over your full dataset multiple times, and being susceptible to the
 264 problems of training neural networks. another solution is favored over this one.

265 It is argued that the translation matrix *must* be an orthogonal matrix [19], to
 266 prove that let's first form the similarity matrix $S = Y(WX^T)$ where Y represents
 267 the word vectors matrix of our source language and X represents the matrix of our
 268 target language, the matrix W is the transformation matrix.

269 Thus the matrix S represents the dot product between every vector in the matrix
 270 Y and its corresponding vector in the transformed matrix WX^T .

271 Looking at the dimensions of the matrices $Y^{n*d} \cdot W^{d*d} \cdot X^{d*n}$ we notice that the
 272 similarity matrix's dimensions are S^{n*n} this represents for each vector in the target
 273 language its similarity with every vector in the source language. We can also have a
 274 reverse similarity matrix between the source vectors and the target vectors defined
 275 by $S' = XQY^T$.

276 For this to remain self consistent, i.e the two similarity matrices have the same
 277 values, because they both calculate the similarity between the two matrices X and
 278 Y , we must have $S' = S^T$ but $S^T = XW^TY^T$ therefore $Q = W^T$ which means that
 279 if W maps the source language to the target language then W^T maps the target
 280 language back to the source language.

281 Given the matrix W we assume that $x \approx W^Ty$ and $y \approx Wx$ therefore $x \approx$
 282 W^TWx therefore we can conclude that the transformation matrix W must be an
 283 orthogonal matrix satisfying $W^TW = I$, where I denotes the identity matrix.

284 We now modify our objective function to satisfy the orthogonality constraint on
 285 W to be:

$$\max_O \sum_{i=1}^n \|y_i^T Wx_i\|^2, \text{ subject to } W^TW = I \quad (7)$$

286 This solution uses Singular Value Decomposition(SVD) to find the translation
 287 matrix. SVD is favored over Principal Components Analysis (PCA) to avoid the
 288 cost of calculating the huge covariance matrix over the dictionaries.

289 The way SVD works is by factoring an $m * n$ matrix M into a product of three
 290 matrices:

$$M = U\Sigma V^T \quad (8)$$

291 Where U is $m * k$, Σ is $k * k$ and V^T is $k * n$. SVD is very common in the field of
 292 natural language processing, it is used in many algorithms including but not limited
 293 to:

- 294 • Latent Semantic Analysis
- 295 • Latent Semantic Indexing
- 296 • Visualization of Word embedding

297 We follow [17]’s method which uses SVD only. Given the source language dictionary
 298 Y_D and the target language dictionary X_D we normalize each vector by dividing it
 299 by its norm. We then multiply the two dictionaries together $M = Y_D^T X_D$ and then
 300 compute the SVD of M which is equal to $M = U \Sigma V^T$.

301 To map the two languages to a single space we multiply each dictionary with
 302 one of the decompositions, our similarity function will become

$$S = YUV^T X^T \quad (9)$$

303 where

$$s_{ij} = y_i^T U V^T x_j \quad (10)$$

304

$$= (U^T y_i) \cdot (V^T x_j) \quad (11)$$

305 thus the two dictionaries are mapped to a single space.

306 This approach produces an exact solution (Equation 9) to our objective function
 307 (Equation 7).

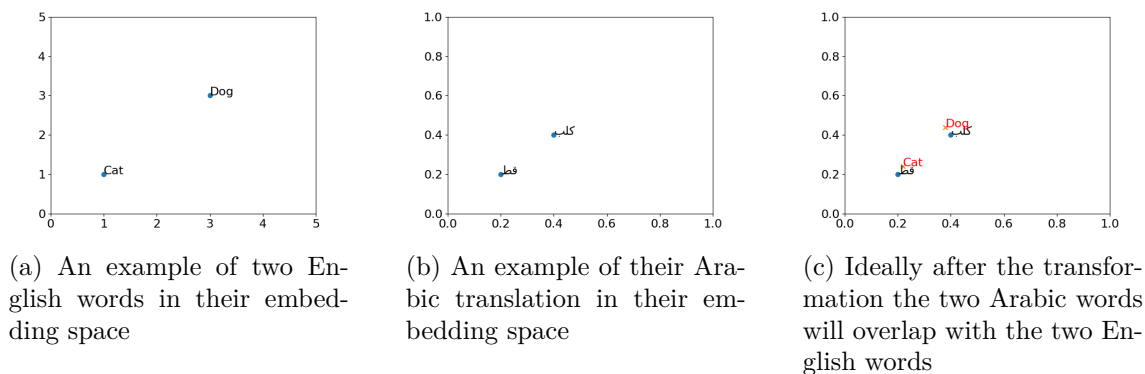


Figure 7: Word embedding transformation

308 4.3. Bi-LSTM encoding

309 Neural Networks have always had a strong ability to perform feature extraction
 310 even with time series or sequences they’re still able to capture time related features
 311 due to their hierarchical nature and recurrence relations.

312 For this we will be using the Long-short term memory [20] algorithm to per-
 313 form feature extraction from our word vectors sequence in order to convey temporal
 314 features to the classifier.

315 LSTM is the core of our model (Figure 8) since it will be responsible for extracting
 316 the useful features from the word vectors to help the classifier find the correct tags
 317 for the sequence.



Figure 8: Model detailed architecture for each time step

318 The main strength of LSTM is its ability to handle long term dependencies in a
 319 sequence of steps while updating its state with information from context due to its
 320 gates that are responsible for updating cell states.

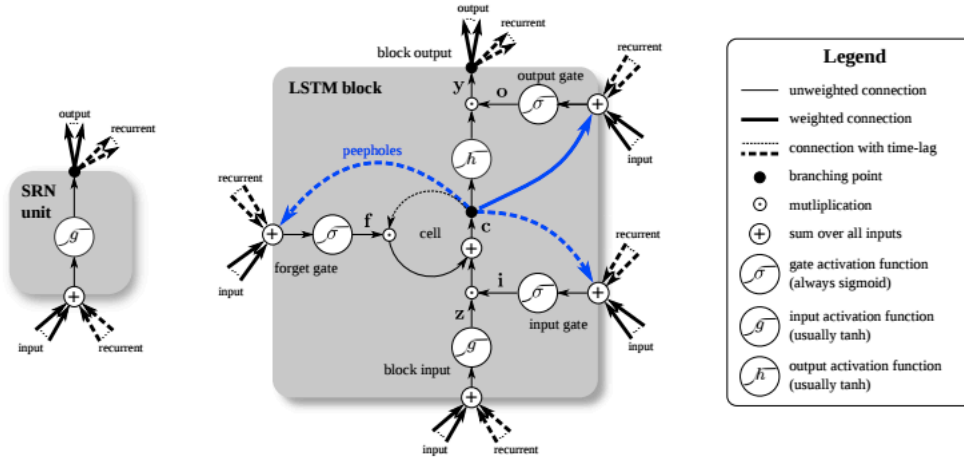


Figure 1. Detailed schematic of the Simple Recurrent Network (SRN) unit (left) and a Long Short-Term Memory block (right) as used in the hidden layers of a recurrent neural network.

Figure 9: LSTM architecture (courtesy of Skyminid AI)

321 Each LSTM unit (Figure 9) holds a cell state C which flows from timestep $t - 1$
 322 to t where some modifications occur in between due to the LSTM gates, these gates
 323 are a way to optionally let information pass through. They are composed of a
 324 sigmoid layer that outputs a value between 0 and 1 and a multiplication operation
 325 in which multiplying by 0 means prevent the information flow and 1 means allow
 326 full information flow. The forget gate receives the hidden state at the previous time
 327 step h_{t-1} and the input x_t and decides how much information needs to be forgotten
 328 by the following equation:

$$f_t = \text{sigmoid}(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (12)$$

329 The input gate decides what new information need to be stored in the cell state, it
 330 follows the following equation:

$$i_t = \text{sigmoid}(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (13)$$

331 We also compute a candidate set of values to be added to the cell state:

$$\tilde{C}_t = \text{tanh}(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (14)$$

332 Following these computations is the update on the cell state which occurs by the
 333 equation:

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (15)$$

334 Finally the model decides on what it should output according to its output gate:

$$o_t = \text{sigmoid}(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (16)$$

335 The output to the hidden state is then computed by:

$$o_t = \text{tanh}(C_t) \quad (17)$$

336 A Bi-directional LSTM layer (Figure 10) is added to our model instead of the
 337 classic Uni-directional LSTM since allowing the model to read text right-to-left and

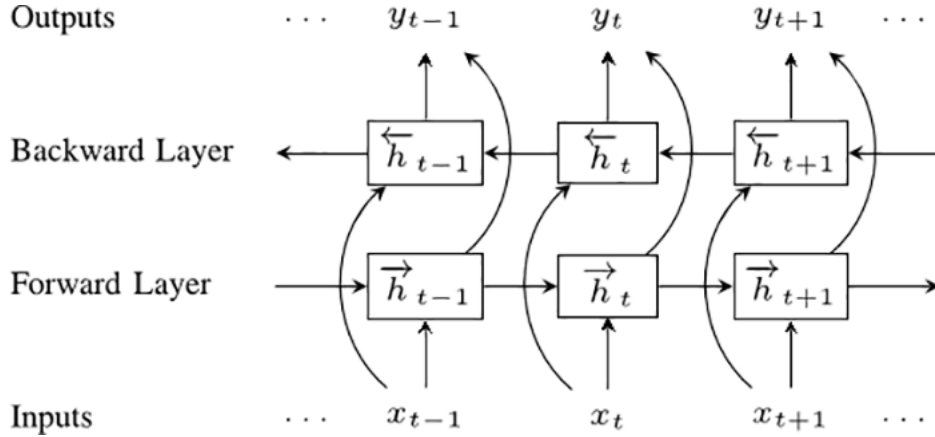


Figure 10: Bi-LSTM architecture [21]

338 left-to-right and aggregating their outputs would improve it significantly by peeking
 339 into the future as well as benefiting from the past [4] [5]. It will also make the model
 340 invariant to language writing direction. for example we train the model on English
 341 which is a left-to-right language and evaluate it on Arabic which is a right-to-left
 342 language.

343 We pass the h_t vector for every time step to the next layer which can be another
 344 layer of Bi-LSTM for an added layer of feature extraction or a Conditional Random
 345 Field classifier.

346 4.4. Conditional Random Field Classifier

347 Going from the features extracted by the Bi-LSTM to class labels requires a
 348 classifier that can take as input the features at each time step and produce as
 349 output the class label for each word.

350 For sequence classification the most common classifiers used in the domain of
 351 NER are:

- 352 • LSTM
- 353 • Hidden Markov Model
- 354 • Simple RNN
- 355 • Conditional Random Field

356 The CRF classifier is used for its strength and its current results in state of the art
 357 NER models [4] [5]. A huge strength of CRF is its ability to work with output classes
 358 that are dependent, an example of this is that for any I-*(Inside class) there must
 359 be a preceding B-*(Beginning class). CRF is capable of exploiting this dependency
 360 thus improving the results and avoiding unintuitive class outputs.

361 Let $h_{1:N}$ be our input to the CRF layer which are the features for each word in
 362 a sequence of length N where (1:N) denotes the words from index 1 to N. These
 363 features are calculated by $h_t = \text{Bi-LSTM}(x_t)$. let $z_{1:N}$ be our class labels for the
 364 input sequence. The linear chain CRF model defines conditional probability of the
 365 class labels given the input by the equation:

$$p(z_{1:N}|h_{1:N}) = \frac{1}{Z} \exp\left(\sum_{n=1}^N \sum_{i=1}^F \lambda_i f_i(z_{n-1}, z_n, h_{1:N}, n)\right) \quad (18)$$

366 In the equation the term Z is the normalization factor or the partition function
 367 which is used to make $p(z_{1:N}|h_{1:N})$ a valid probability, the term F is the number of
 368 feature functions we're going to use.

$$Z = \sum_{z_{1:N}} \exp\left(\sum_{n=1}^N \sum_{i=1}^F \lambda_i f_i(z_{n-1}, z_n, h_{1:N}, n)\right) \quad (19)$$

369 In the probability equation we had feature functions f_i , each of which is a function
 370 that takes as input the current class label, the label of the previous class, the input
 371 and the current position and outputs a boolean value 0 or 1. λ_i is the feature weight
 372 for feature i , if this weight is large and positive then we emphasize that the probability
 373 of the current label is high if this feature is true. An example of a feature function:

$$374 \quad f_1(z_{n-1}, z_n, h_{1:N}, n) = \begin{cases} 1, & \text{if } z_{n-1} = \text{B-PERS and } h_t = \text{Bi-LSTM(John)} \\ 0, & \text{Otherwise} \end{cases} \quad \text{A posi-}$$

375 tive weight for this function would indicate that the model prefers the tag B-PERS
 376 for the current word, if the function was assigned a negative weight it would indicate
 377 that the model does not prefer the tag B-PERS for the current word.

378 In our case the features will be the output of the Bi-LSTM. In this case the
 379 Bi-LSTM layer(s) can be considered as the feature function we are using.

380 The appropriate objective for parameter learning of CRF is to maximize the con-
 381 ditional likelihood of the training data where m is the number of training examples:

$$\sum_{j=1}^m \log p(z_{1:N}|h_{1:N}) \quad (20)$$

382 This learning procedure is conducted by computing the gradients for the objective
 383 function (Equation 20) and use the gradient in a gradient based optimization algo-
 384 rithm like Stochastic Gradient Descent.

385 Now that the model is trained we can use it to calculate the probability $p(z_{1:N}|h_{1:N})$
 386 for any tag for each token in an input sequence. An approach to take the final se-
 387 quence is by greedily taking the tag that has the highest probability. However this
 388 approach would have to check an exponential number of tags because it has to check
 389 K^N possible tags with K being the number of classes in our model. Another ap-
 390 proach would be to use the Viterbi algorithm [22] which is a dynamic programming
 391 algorithm that works in polynomial time to find the optimal sequence of tags.

392 5. Experimental Results

393 First we define the metrics we used in our evaluation. We used Precision and
 394 Recall and their harmonic mean (F1-score).

395 Precision is defined as the number of true positives divided by the number of
 396 true positives and false positives. It is considered as a measure of the classifier's
 397 exactness.

$$Precision = \frac{TruePositives}{TruePositives + FalsePositives} \quad (21)$$

398 Recall is defined as the number of true positives divided by the number of true
 399 positives and false negatives. It can be considered the sensitivity or the true positive
 400 rate.

$$Recall = \frac{TruePositives}{TruePositives + FalseNegatives} \quad (22)$$

401 In our evaluation we use the F1-score which is a balance between them.

$$F1score = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (23)$$

402 We trained our model on the WikiNER English dataset [18] by using the word
 403 embedding model [16] which is open source. We evaluated the model on English
 404 and noted its results then we tested it on the ANER dataset [3] with only applying
 405 the transformation and no fine tuning of the model.

406 The model was written using Keras [23], it was trained on an Ubuntu 16.04
 407 system running on an Intel i7-7500U processor and accelerated by an Nvidia Geforce
 408 940MX GPU. All the code for the experiment is available as open source on Github:
 409 [Code on Github](#)

410 The results of the model imply that the model succeeded at detecting many
 411 named entities in the Arabic text without being exposed to Arabic before. This
 412 approach is valid for any other language that we have a transformation matrix for.
 413 Recently 78 matrices for transformation [17] have been released as open source.
 414 These matrices can be used to transform the word embedding space from any of the
 415 78 languages to English, which can be used to detect named entities in 78 different
 416 languages with only being trained on English.

Classification results per class for English			
	Precision	Recall	F1-Score
PER	0.88	0.93	0.91
MISC	0.80	0.73	0.76
LOC	0.86	0.85	0.86
ORG	0.82	0.72	0.76
Average scores	0.84	0.82	0.83

Table 1: Classification Report for Precision, Recall and F1-Scores on WikiNER

417 We can see from Table 1 that the model successfully learns to detect named
 418 entities in the English text. The model’s accuracy depends on many factors among
 419 which the training data and the complexity of the model. The model seems to
 420 perform best on the persons class and seems to perform worst on the miscellaneous
 421 class.

Classification results per class for Arabic			
	Precision	Recall	F1-Score
PER	0.05	0.01	0.01
MISC	0.01	0.57	0.02
LOC	0.07	0.00	0.00
ORG	0.02	0.05	0.03
Average scores	0.05	0.07	0.02

Table 2: Classification Report for Precision, Recall and F1-Scores on ANER without Alignment

422 In Table 2 the model is tested on ANER which is the Arabic dataset without
 423 aligning the word vectors with the English word vectors. Recall that the model
 424 was never trained to understand Arabic and so testing the model on Arabic would

425 produce a score of 0 (or almost 0) which is the case as seen in the table. The model
 426 fails to capture any entities -Except for a very tiny number of entities that could be
 427 random-.

Classification results per class for Arabic			
	Precision	Recall	F1-Score
PER	0.68	0.52	0.59
MISC	0.05	0.06	0.06
LOC	0.80	0.44	0.57
ORG	0.32	0.08	0.12
Average scores	0.58	0.36	0.43

Table 3: Classification Report for Precision, Recall and F1-Scores on ANER with Alignment

428 The results in Table 3 shows that after alignment the model succeeds at detecting
 429 named entities in Arabic which it was never trained on. It detects the class persons
 430 with precision of 68% compared to 88% in English and it's also its top scoring class.

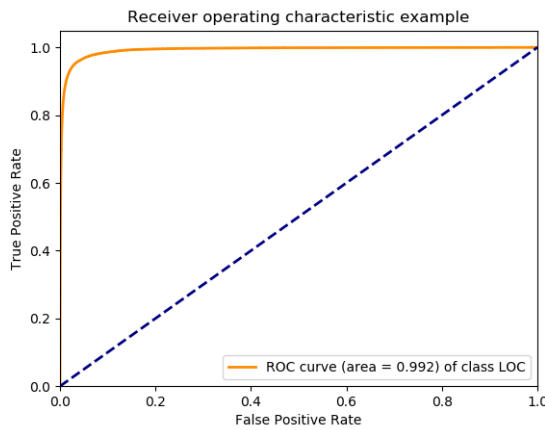


Figure 11: ROC curve for LOC class in English

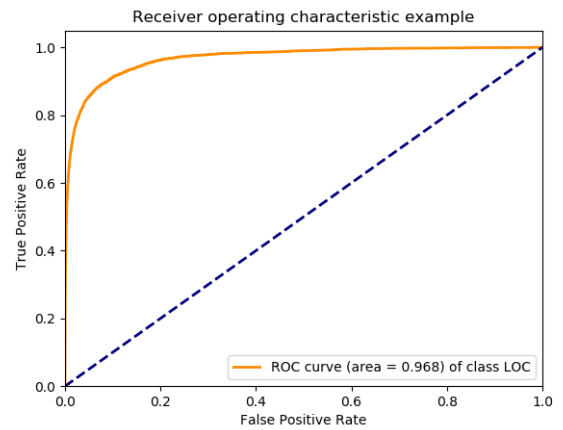


Figure 12: ROC curve for LOC class in Arabic

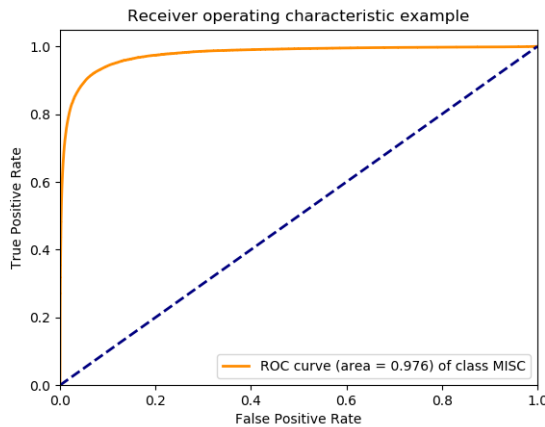


Figure 13: ROC curve for MISC class in English

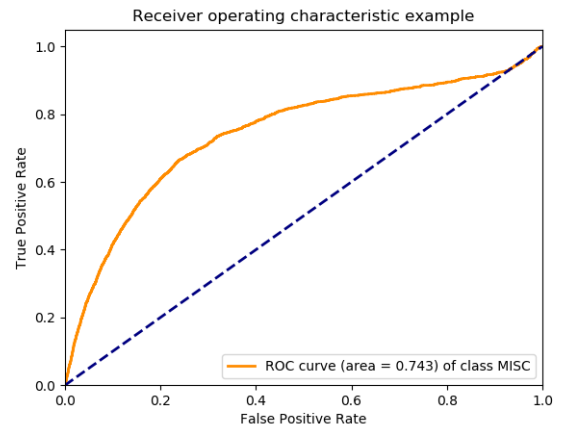


Figure 14: ROC curve for MISC class in Arabic

431 ROC figures and their corresponding AUC in Figures 11, 12, 13, 14, 15, 16,
 432 17, 18 show that the proposed approach performs very well compared to random

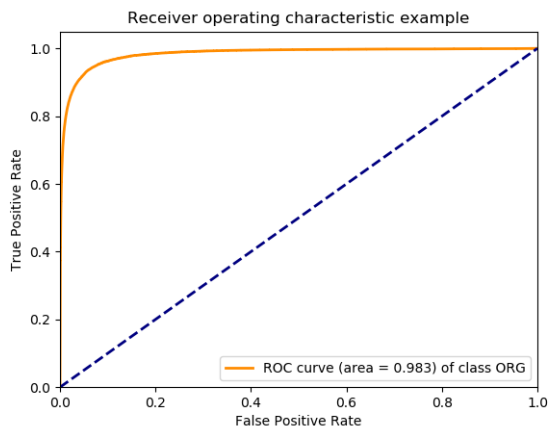


Figure 15: ROC curve for ORG class in English

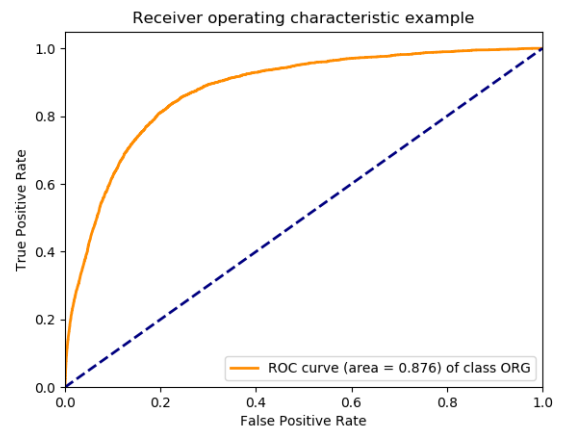


Figure 16: ROC curve for ORG class in Arabic

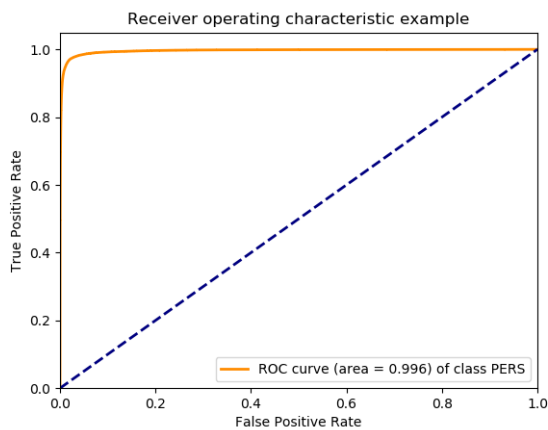


Figure 17: ROC curve for PERS class in English

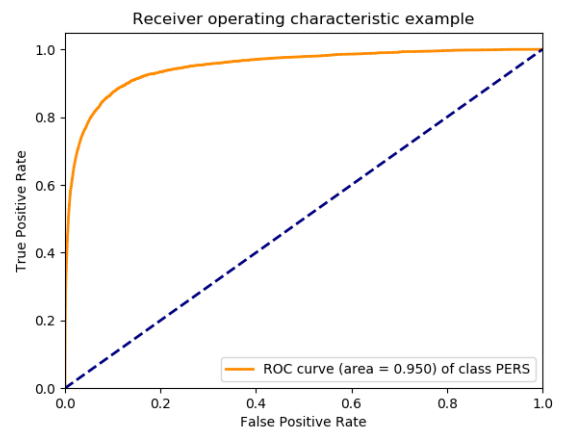


Figure 18: ROC curve for PERS class in Arabic

433 guessing and that the system is capable of classifying the aforementioned classes in
 434 both English and Arabic without the need for retraining the model on Arabic.

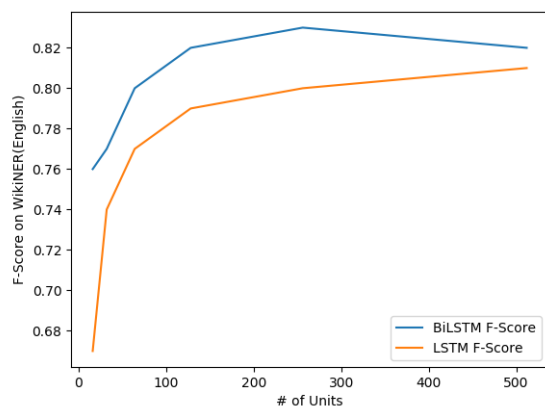


Figure 19: LSTM vs BiLSTM performance on English data

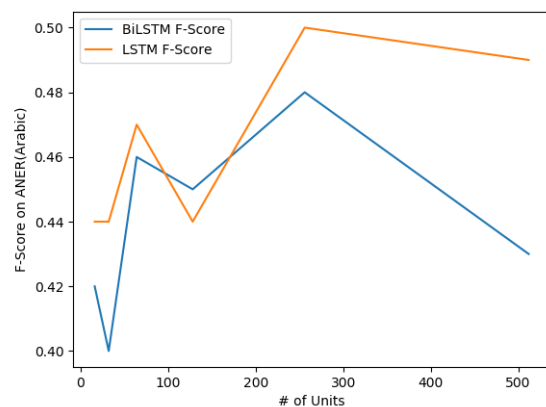


Figure 20: LSTM vs BiLSTM performance on aligned Arabic data

435 Hyper parameter selection for the model was performed empirically by iterating

436 over a finite set of parameters and choosing values that maximize our objective
437 function (Grid search).

438 The main set of parameters to choose from were whether to use Bidirectional
439 LSTM or Unidirectional LSTM and the number of units at the LSTM layer. These
440 parameters were tested on both the English and Arabic models.

441 Figure 19 shows that BiLSTM performance is superior to LSTM performance
442 on English data while it fluctuates on Arabic data according to Figure 20. It's also
443 noted that the F-score values starts degrading just below 300 units. This makes the
444 optimal choice for units to be around 256 units.

445 6. Conclusion

446 We exploit orthogonal transformation of word embeddings to create a language
447 independent NER model that was trained on English and evaluated on Arabic with-
448 out being explicitly exposed to Arabic before or fine tuned on an Arabic dataset.
449 This paves the road to language independent NLP models that are capable of solving
450 typical NLP tasks without being explicitly trained on a specific language.

451 References

- 452 [1] Jiafeng Guo, Gu Xu, Xueqi Cheng, and Hang Li. Named entity recognition
453 in query. In *Proceedings of the 32nd international ACM SIGIR conference on*
454 *Research and development in information retrieval*, pages 267–274. ACM, 2009.
- 455 [2] Diego Mollá, Menno Van Zaanen, Steve Cassidy, et al. Named entity recognition
456 in question answering of speech data. 2007.
- 457 [3] Yassine Benajiba, Paolo Rosso, and José Benedíruiz. Anersys: An arabic named
458 entity recognition system based on maximum entropy. *Computational Linguistics*
459 *and Intelligent Text Processing*, pages 143–153, 2007.
- 460 [4] Zhiheng Huang, Wei Xu, and Kai Yu. Bidirectional lstm-crf models for sequence
461 tagging. *arXiv preprint arXiv:1508.01991*, 2015.
- 462 [5] Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya
463 Kawakami, and Chris Dyer. Neural architectures for named entity recogni-
464 tion. *arXiv preprint arXiv:1603.01360*, 2016.
- 465 [6] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recog-
466 nition with deep recurrent neural networks. In *Acoustics, speech and signal*
467 *processing (icassp), 2013 ieee international conference on*, pages 6645–6649.
468 IEEE, 2013.
- 469 [7] John Lafferty, Andrew McCallum, and Fernando CN Pereira. Conditional ran-
470 dom fields: Probabilistic models for segmenting and labeling sequence data.
471 2001.
- 472 [8] Lance A. Ramshaw and Mitchell P. Marcus. Text chunking using
473 transformation-based learning. *CoRR*, cmp-lg/9505040, 1995.
- 474 [9] Michal Ziemski, Marcin Junczys-Dowmunt, and Bruno Pouliquen. The united
475 nations parallel corpus v1. 0. In *LREC*, 2016.

- 476 [10] Kareem Darwish. Named entity recognition using cross-lingual resources: Ara-
477 bic as an example. In *ACL (1)*, pages 1558–1567, 2013.
- 478 [11] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas,
479 Pablo Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleeef, Sören
480 Auer, and Chris Bizer. DBpedia - a large-scale, multilingual knowledge base
481 extracted from wikipedia. *Semantic Web Journal*, 2014.
- 482 [12] Zhilin Yang, Ruslan Salakhutdinov, and William Cohen. Multi-task cross-
483 lingual sequence tagging from scratch. *arXiv preprint arXiv:1603.06270*, 2016.
- 484 [13] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Em-
485 pirical evaluation of gated recurrent neural networks on sequence modeling.
486 *arXiv preprint arXiv:1412.3555*, 2014.
- 487 [14] Zellig S Harris. Distributional structure. *Word*, 10(2-3):146–162, 1954.
- 488 [15] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean.
489 Distributed representations of words and phrases and their compositionality.
490 In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- 491 [16] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enrich-
492 ing word vectors with subword information. *arXiv preprint arXiv:1607.04606*,
493 2016.
- 494 [17] Samuel L Smith, David HP Turban, Steven Hamblin, and Nils Y Hammerla.
495 Offline bilingual word vectors, orthogonal transformations and the inverted
496 softmax. *arXiv preprint arXiv:1702.03859*, 2017.
- 497 [18] Joel Nothman, Nicky Ringland, Will Radford, Tara Murphy, and James R.
498 Curran. Learning multilingual named entity recognition from Wikipedia. *Ar-*
499 *tificial Intelligence*, 194:151–175, 2012.
- 500 [19] Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. Normalized word embedding
501 and orthogonal transform for bilingual word translation. In *HLT-NAACL*, pages
502 1006–1011, 2015.
- 503 [20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural*
504 *computation*, 9(8):1735–1780, 1997.
- 505 [21] Alex Graves, Navdeep Jaitly, and Abdel rahman Mohamed. Hybrid speech
506 recognition with deep bidirectional lstm. *2013 IEEE Workshop on Automatic*
507 *Speech Recognition and Understanding*, pages 273–278, 2013.
- 508 [22] G David Forney. The viterbi algorithm. *Proceedings of the IEEE*, 61(3):268–278,
509 1973.
- 510 [23] François Chollet et al. Keras. <https://github.com/fchollet/keras>, 2015.