# Improving Neural Sequence Labelling using Additional Linguistic Information

Mahtab Ahmed, Muhammad Rifayat Samee, Robert E. Mercer

*Department of Computer Science*
*University of Western Ontario*
London, Ontario, Canada
mahme255, msamee, rmercer@uwo.ca

*Abstract*—Sequence labelling is the task of assigning categorical labels to a data sequence. In Natural Language Processing, sequence labelling can be applied to various fundamental problems, such as Part of Speech (POS) tagging, Named Entity Recognition (NER), and Chunking. In this study, we propose a method to add various linguistic features to the neural sequence framework to improve sequence labelling. Besides word level knowledge, sense embeddings are added to provide semantic information. Additionally, selective readings of character embeddings are added to capture contextual as well as morphological features for each word in a sentence. Compared to previous methods, these added linguistic features allow us to design a more concise model and perform more efficient training. Our proposed architecture achieves state of the art results on the benchmark datasets of POS, NER, and chunking. Moreover, the convergence rate of our model is significantly better than the previous state of the art models.

*Index Terms*—Sequence Labelling, Long Short Term Memory, Conditional Random Field, Linguistic Features.

## I. Introduction

Linguistic sequence labelling is one of the first tasks focusing on natural language processing using deep learning and it has been well examined over the past decade [1]–[3]. Part of speech (POS) tagging, named entity recognition (NER), and chunking are subclasses of sequence labelling. They play a vital role in fulfilling many downstream applications, such as relation extraction, syntactic parsing, and entity linking [4]–[6]. POS tagging assigns a tag to each word in a text, where a tag represents the lexical category of a word. NER is a subtask of information extraction that seeks to locate and classify named entities in text. Chunking identifies the POS and short phrases in a sentence by doing shallow parsing and also groups words into syntactically correlated phrases. These labelled texts can later be used for different applications such as machine translation, information retrieval, word sense disambiguation, and natural language understanding etc.

Before neural sequence models, traditional algorithms were based on Hidden Markov Models (HMMs) [7], [8] and Conditional Random Fields (CRFs) [9], [10]. The problem with these models is they are heavily dependent on manually handcrafted features. So it becomes difficult to apply them in real life applications because it is not practically possible to always have human expertise.

To overcome these drawbacks, Neural Network (NN) based models have been proposed in which the models are responsible for extracting higher level features from the data [11], [12]. Recurrent Neural Networks (RNNs) along with its variants, Long Short Term Memory (LSTM) [13], [14] and Gated Recurrent Units (GRUs) are found to work very well with sequence data as they can capture long distance dependencies [14]–[16]. Nevertheless, considering the overwhelming number of their parameters and the relatively small size of most human annotated sequence labelling corpora, annotations alone may not be sufficient to train complicated models. So, guiding the learning process with extra knowledge could be a wise choice [17], [18]. For example, before tagging the word '*flies*' as either a verb or a noun in the sentence '*Time flies like an arrow*', having its semantic meaning would make a correct tagging straightforward.

Knowing the sense of a word prior to sequence labelling (POS or NER) often gives the most probable tag for that word. Word senses can be obtained from a variety of sources: WordNet [19], a lexical database for English that can be queried for the sense of a word given its context; the simplified LESK algorithm [20], [21] which uses the dictionary definition of each word in a sentence as extra context to suggest the word sense; and linear algebraic methods, one [22] which uses a random walk on a discourse model and represents the vector of the base word as a linear combination of its probable sense vectors.

In this paper, we propose a novel deep neural architecture for doing sequence labelling incorporating not only semantic features through word senses but also the rich morphology of the words. We provide an in depth analysis of the design of this architecture giving some insights regarding how each feature is introduced into the architecture. Our sequence model achieves state of the art results on the three sequence labelling tasks, POS, NER, and chunking, and has a training time at least four times faster than the currently available state of the art models.

## II. RELATED WORK

Huang et al. [23] propose a few models for the sequence tagging task. Apart from just word embeddings, they first use morphological as well as bigram and trigram information as their input features. Later, they use LSTM and Bidirectional LSTM (BLSTM) with CRF to do the final tagging. Lample et al. [24] extract character embeddings from both the left and right directions, concatenate these with word embeddings and then use a stacked LSTM along with CRF to do the tagging. Liu et al. [25] propose a model leveraging both word as well as character level features. It includes a language model to represent the character level knowledge along with a highway layer to avoid the feature collision. Finally it is trained jointly as a multitask learning.

Yu et al. [26] propose a general purpose tagger using a convolutional neural network (CNN). First, they use CNNs to extract the character level features and then concatenate it with word embeddings, position embeddings and binary features. Finally they use another CNN to get the contextual features as well as to do the tagging. Ma and Hovy [27] propose an end-to-end sequence labelling model using a combination of BLSTM, CNN and CRF. They use a CNN to get the character level information, concatenate it with word embeddings and then apply BLSTM to model the contextual information. Finally they generate the tags by using a sequential CRF layer. Rei [28] trains a language model type objective function using BLSTM-CRF to predict the surrounding words for every word in the corpus and utilizes it for sequence labelling.

The contribution of this paper combines the common themes found in these previous works (morphology encoded as character embeddings, and word embeddings) with word senses in a new architecture that integrates these embeddings and the outputs of a CNN, a BLSTM, and a CRF in novel ways.

## III. THE MODEL: BLSTM-CRF

In this section, we describe our work in detail. We first explain each of the pieces of the complete architecture and then we explain how we combine those pieces to build our model. This section also explains the morphological and semantic features that we have added with our model to get the improved performance that is discussed in Section V.

### A. Recurrent unit: Bidirectional LSTM

Recurrent neural networks (RNNs) are the best known and most widely used NN model for sequence data as they go over the entire sequence through time and try to remember it in a compressed form. Although its variant, LSTM, is very good with long term dependencies, for many sequence labelling task, it is important to keep track of these dependencies from the future as well as from the past. But LSTM has just one hidden state from the past and changes that hidden state recursively through time. An elegant solution to this problem is going over the sequence in both forward and backward directions with two hidden states and finally concatenating the output from both directions. This bidirectionality has proven to be very effective in some prior works [29]–[31]. The

resulting network, the Bidirectional LSTM (BLSTM), is the RNN variant that is used by the model described below.

### B. Word Sense

Knowing the sense of a word prior to tagging makes the tagging task more straightforward. Generally, polysemy is captured in standard word vectors, but the senses are not represented as multiple vectors. So we have trained an adaptive skip gram model, AdaGram, [32] which gives a vector for each sense of a word. It is a non-parametric Bayesian extension of the skip-gram model and is based on the constructive definition of Dirichlet process (DP) [33]. It can learn the required number of representations of a word automatically.

In our model, we denote a set of input words as $X = \{x_i\}_{i=1}^N$ and their context as $Y = \{y_i\}_{i=1}^N$. The $i$th training pair $(x_i, y_i)$ consists of words $x_i = o_i$ with context $y_i = (o_t)_{t \in c(i)}$, where $c(i)$ is the index of the context words. Then, instead of maximizing the probability of generating a word given its contexts [34], we maximize the probability of generating the context given its corresponding input words [32]. Our final objective function becomes,

$$p(Y|X,\theta) = \prod_{i=1}^N p(y_i|x_i,\theta) = \prod_{i=1}^N \prod_{j=1}^C p(y_{ij}|x_i,\theta) \quad (1)$$

where, $\theta$ is the set of model parameters. The drawbacks of this objective function is that it captures just one representation of a word which goes against a word having different senses depending on the context [22]. To counter this, AdaGram introduces a new latent variable $z$ which captures the required number of senses even though the number of structure components of the data is unknown a priori. In AdaGram, if the similarities of a word vector with all its existing sense vectors are below a certain threshold, a new sense is assigned to that word with a prior probability $p$. The prior probability of the $k$th meaning of word $w$ is

$$p(z = k|w, \beta) = \beta_{wk} \prod_{r=1}^{k-1} (1 - \beta_{wr}),$$

$$p(\beta_{wk}|\alpha) = Beta(\beta_{wk}1, \alpha), k = 1 \dots \quad (2)$$

where $\beta$ is a latent variable and $\alpha$ controls the number of senses. Theoretically, it is possible to have an infinite number of senses for each word $w$. However, as long as we have a finite amount of data, the number of senses can not be more than the number of occurrences of that word. With more data, it can increase the complexity of the latent variables thereby allowing more distinctive meanings to be captured. Taking all the facts into account, our final objective function becomes,

$$p(Y, Z, \beta|X, \alpha, \theta) = \prod_{w=1}^V \prod_{k=1}^\infty p(\beta_{wk}|\alpha)$$

$$\prod_{i=1}^N [p(z_i|x_i,\beta) \prod_{j=1}^C [p(y_{ij}|z_i, x_i, \theta)] \quad (3)$$

where $Z = \{z_i\}_{i=1}^N$ is a set of senses for all the words.

## C. Convolutional Neural Network

Convolutional Neural Networks (CNNs) are good for extracting n-gram features from a sentence [35]. They consist of kernels (i.e., a matrix of weights) which are used to go over the input word embedding matrix with a variable stride length and extract some higher level features. In our model, we use a CNN to get the bigram features. To do that, first we pad the input sentence [1] and pass it to an embedding layer. This layer represents the sentence as a matrix of size $(m+1) \times d$, where $m$ is the actual sentence length and $d$ is the embedding dimension. Next, we initialize a kernel of size $2 \times d$ with stride length 1 and convolve it with the input sentence matrix. This results in a bigram embedding matrix $B$ of size $m \times d$ using Eqn. 4.

$$B_{i,:} = \sum_{j=1}^{2} I_{i+j,:} * K_{j,:} \tag{4}$$

where $I$ is the input sentence matrix, $K$ is the convolution kernel and $n$ is the maximum sequence length for the current batch. Later, this bigram embedding is passed to a BLSTM layer to extract more abstract features.

## D. Conditional Random Field

Each of the tasks that we are modelling requires a tag to be assigned to each word. In addition to using the current word to predict its tag, it is also possible to use the information about the neighboring words' tags. There are two main ways to do this. One way is to calculate the distribution of tags over each time step and then use a beam search-like algorithm, such as maximum entropy markov models [36] and maximum entropy classifiers [37], to find the optimal sequence. Another way is to focus on the entire sentence rather than just the specific positions which leads to Conditional Random Fields (CRFs) [9]. CRFs have proven to give a higher tagging accuracy in cases where there are dependencies between the labels. Like the bidirectionality of BLSTM networks a CRF can provide tagging information by looking at its input features bidirectionally.

In our model we denote a generic input sequence as $x = \{x_i\}_{i=1}^{N}$, generic tag sequence as $y = \{y_i\}_{i=1}^{N}$, and set of possible tag sequences of $x$ as $F(x)$. Then we use CRF to calculate the conditional probability over all possible tag sequences $y$ given $x$ as

$$p(y|x; W, b) = \frac{\prod_{i=1}^{n} \phi_i(y_{i-1}, y_i; x)}{\sum_{y' \in F(x)} \prod_{i=1}^{n} \phi_i(y'_{i-1}, y'_i; x)} \tag{5}$$

where $\phi(.)$ is the score function for the transition between the tag pair $(y', y)$ given $x$. We train this CRF model using maximum likelihood estimation (MLE) [38]. For a training pair $(x_i, y_i)$ we maximize

$$L(W, b) = \sum_{i} \log p(y|x; W, b) \tag{6}$$

[1] The '*valid*' convolution operation reduces the size of the feature matrix by $k - L$, where $k$ is the size of the kernel vector and $L$ is the stride length. For us, $L = 1$. To keep the size of the feature matrix uniform through the model, we padded a start token $<start>$ at the beginning of the input sentence.

where $W$ is the weight matrix and $b$ is the bias term. While decoding, we search for the best tag $\hat{y}$ with the highest conditional probability using the Viterbi algorithm [39].

$$\hat{y} = \arg\max_{y \in F(x)} p(y|x; W, b) \tag{7}$$

## E. Morphology: Spelling and suffix features

For the morphological features, we have focused on spelling features and suffix features.

We extract the following 14 spelling features for a given word and store it as a binary vector $SV_{1 \times 14}$.

- whether it is composed only of alphabetics or not
- whether it contains non-alphabetic characters except '.' or not
- whether it starts with a capital letter or not
- whether it is composed only of upper case letters or not
- whether it is composed only of lower case letters or not
- whether it is composed only of digits or not
- whether it is composed of a mixture of alphabetics and numbers or not
- whether it is the starting word in the sentence or not
- whether it is the last word in the sentence or not
- whether it is in the middle of the sentence or not
- whether it ends with an apostrophe s ('s) or not
- whether it has punctuation or not
- whether it is the first word in the sentence and starts with a capital letter or not
- whether it is composed mostly of digits or not

Apart from extracting these features, we also replace all the numbers in the corpus with the $<number>$ tag.

We have assembled a list of 137 suffixes from https://www.learnthat.org/pages/view/suffix.html and have used the ten that occur most often in our corpus for this study. Then for each of these suffixes, we have collected the words that end with that suffix and have recorded their POSs as well as the frequency. Next, we made an assumption that if a word $w$ with POS $x$ ends with a specific suffix $s$ exceeds a frequency threshold in the training set, then $s$ is the true suffix of word $w$. We record the pair as $(w, s)$. Finally, we create a one hot vector $SUV_{1 \times 10}$ for each word where a 1 at index $k$ means the word has the $k$th suffix.

## F. BLSTM-CRF model

In this sub-section, we combine the BLSTM and CRF models with some feature connection techniques to form our final BLSTM-CRF model. We divided this final model into some modules and the description of each of these modules is as follows:

**Module 1: Word Level Features** This module starts with an embedding layer. In detail, we initialize the emebedding layer randomly as well as using pre-trained embeddings (GloVe / word2vec). Next we represent each sentence as a column vector $I_{m \times 1}$ where each element of the vector is a unique index of the corresponding word. Then we pass

this vector to an embedding layer which gives a matrix representation $W_{m \times d}$. Here, $d$ is the embedding dimension.

**Module 2: Character Embedding** In this module, first we split a word into its characters and then transform it into a column vector $C_{k \times 1}$, where $k$ is the word length and each element of the vector is a unique index of the corresponding character. Next we initialize an embedding layer randomly and pass the character vector into it. This will change the representation to a matrix of size $k \times n$ where $n$ represents the embedding dimension. Then we use an LSTM on this matrix and store the last hidden state of this LSTM as the character level representation $C_{1 \times n}$ of the word. Finally, for a sentence with $m$ words, it is stored as a matrix $C_{m \times n}$.

**Module 3: Selective Pickup from Char LSTM (SP-CLSTM)** In this module, we introduce a new way of capturing the morphological features as well as the context features. The word embeddings from module 1 gives the contextual features in both directions and the character embeddings from module 2 gives the lexical information. We capture both sets of information by first representing each sentence in terms of its characters $I_{(k \times m) \times 1}$ and then turn this into a matrix of size $k \times m \times d$ through a random embedding layer. Then we apply a BLSTM over this representation and finally we pick those indices from the output where each word ends. This selective pickup provides the morphological information of a word as well as information about the previous words in the sequence.

$$\tilde{C}_{m \times d} = \text{SELECT}(\text{BLSTM}(I_{(k \times m) \times d})) \quad (8)$$

**Module 4: Sense Features** This module calculates the sense level contextual features of a sentence. First, we initialize a sense embedding layer using the pre-trained sense embeddings from AdaGram. Then we tag each word in the input sentence using the module `disambiguate` from AdaGram (the word '*apple*' with sense 2 is tagged as '*apple_2*'). This modified input sentence is then passed to the embedding layer initialized before and finally the resultant output is passed to a BLSTM layer. The output of this BLSTM layer gives the sense level contextual feature $S_{m \times d}$.

**Module 5: Bigram Features** This module calculates the bigram embedding features $B_{m \times d}$ of a sentence as described in the Subsection III-C.

**Module 6: The Connection Technique** In this module, we combine all the features and the modules using some novel connection techniques and build our final BLSTM-CRF model as shown in Figure 1. First we concatenate the word embedding from module 1 with the character embedding from module 2 and the suffix vector from Subsection III-E as $[W_{m \times d}, C_{m \times n}, SUV_{m \times 10}]$. Following this, we apply a BLSTM on this new embedding matrix, calling this output $O^1_{m \times d}$. The outputs of modules 3, 4 and 5 are called $O^2_{m \times d}$, $O^3_{m \times d}$, and $O^4_{m \times d}$, respectively. Then we initialize four scalar weights $w_1$, $w_2$, $w_3$ and $w_4$ with initial value 1.0 and add them as model parameters. We form a linear combination of the $w_i$ weighted $O^i$s to form the final output.

$$O = \sum_{i=1}^{4} O^i w_i \quad (9)$$

The final output ($O_{m \times d}$) have pieces of information from all the features that we calculated above. We choose linear addition rather than concatenation of these output features, because concatenation will result in a very large feature matrix and the network have to tune each of the cell of this matrix during back-propagation. Following this, we initialize an LSTM layer where we pass the final output from Eqn. 9 at each time step $\tilde{O}^i_{1 \times z} = \text{LSTM}(O^i_{1 \times d}, h^{i-1}_{1 \times d})$ and store the outputs separately $\tilde{O}_{m \times z} = [\tilde{O}^1, \tilde{O}^2, \ldots, \tilde{O}^m]$. This LSTM layer unfolds at each time step taking the hidden state of the previous time step to initialize the hidden state of the current time step. The previous hidden state has the information about the previous tag and initializing the current hidden state with the previous one explicitly gives this information. Next we pass the output from each time step to a `tanh` layer $T_{1 \times d} = \text{tanh}(\tilde{O}^i)$, which squeezes the values between $[-1, 1]$. Then we concatenate this `tanh` output $T_{1 \times d}$ with the spelling features $SF_{1 \times 14}$ calculated in subsection III-E and pass this to a fully connected (FC) layer. This FC layer maps the output to the number of tag classes $Y_{1 \times c} = \text{FC}([T_{1 \times d}, SUV_{1 \times 14}])$, where $c$ represents number of classes. We do this for each time step and concatenate the results to make a final tensor $Y_{m \times c}$. Finally, we pass this tensor to the CRF layer and calculate the possible tag sequence for the given input sequence.

## IV. EXPERIMENTAL SETUP

In this section, we describe the detailed experimental setup for the evaluation of our study. We first explain the dataset statistics for each tagging task. Following this, we explain the working environment details along with the hyper-parameter settings of our architecture.

### A. Dataset Description

We test our BLSTM-CRF model on three NLP tagging tasks: Penn TreeBank (PTB) POS tagging, CoNLL 2000 chunking, and CoNLL 2003 NER. Table I shows the number of sentences in the training, validation and test sets respectively for each corpus. We utilize the BIO2 explanation standard for the chunking and NER tasks.

|       | WSJ   | CoNLL00 | CoNLL03 |
|-------|-------|---------|---------|
| Train | 39831 | 8936    | 14987   |
| Valid | 1699  | N/A     | 3466    |
| Test  | 2415  | 2012    | 3684    |

TABLE I: Dataset Description

Table II shows the detailed hyper-parameter settings of our model and some of the hyper-parameters for AdaGram (the remaining parameters are set to their default values [40]). We train our model on Nvidia GeForce GTX 1080 GPU with both the 'Adam' and 'SGD' optimizers. All of the results in the next section are reported using 'SGD' as it was giving the best results. The 'Learning rate decay' parameter was only used with the 'SGD' optimizer. We used PyTorch 0.3.1 to

Fig. 1: BLSTM-CRF model architecture

| BLSTM-CRF | |
|---|---|
| **Hyper-parameter** | **Range Selected** |
| Learning rate | 0.001 / 0.015 / 0.01 |
| Batch size | 10 / 50 / 100 |
| No. of LSTM layers | 1 / 2 / 3 |
| Momentum | 0.9 |
| Dropout | 0.5 / 0.2 / 0.1 |
| Word embedding size | 300 / 200 / 100 |
| Character embedding size | 50 / 30 |
| Initial scalar weight value | 1.0 |
| Gradient clipping | 5 / 20 / 50 |
| Weight decay | $10^{-5}$ |
| Learning rate decay | 0.05 |
| CNN kernel size | $2 \times (300/200/100)$ |
| **AdaGram** | |
| Epoch | 1000 |
| Window size | 10 / 7 / 5 |
| No. of prototypes | 5 |
| Sense embedding size | 300 |
| Prior prob. of new sense | 0.1 |
| Initial weight on first sense | -1 |
| Word embedding size | 300 / 200 / 100 |

TABLE II: Ranges of different hyper-parameters searched during tuning.

implement our model and Julia 0.4.5 for running AdaGram under the Linux environment.

## V. EXPERIMENTAL RESULTS

In this section, we describe in detail the results obtained with our proposed architecture. As the evaluation metrics, we use accuracy for the WSJ corpus and F1 score (micro averaged) for the CoNLL00 and CoNLL03 tasks. This section also contains the results of the top performing models for all three sequence labelling tasks. Additionally, we show the rate of convergence of our model compared to the state of the art one. Finally, we conclude this section by giving an ablation study by removing certain modules as well as features and mixing them in different combinations.

Table III shows the performances of all of the chunking systems. An SVM based classifier [41] won the CoNLL 2000 challenge with an F1 score of 93.48%. However, later they improved their result up to 93.91% [42]. Recently, most of the models incorporate CRF in their architecture to capture the tag

| Model | F1-score |
|---|---|
| SVM classifier [41] | 93.48 |
| SVM classifier [42] | 93.91 |
| BI-LSTM-CRF [23] | 94.13 |
| Second order CRF [43] | 94.29 |
| Second order CRF [39] | 94.30 |
| Conv. network tagger [44] | 94.32 |
| Second order CRF [45] | 94.34 |
| BLSTM-CRF (Senna) [23] | 94.46 |
| HMM + voting [46] | 95.23 |
| BLSTM-CRF (Ours) | **96.76** |

TABLE III: Comparison of F1 scores of different models for chunking

| Model | F1-score |
|---|---|
| Conv-CRF [44] | 81.47 |
| BLSTM-CRF [23] | 84.26 |
| MaxEnt classifier [47] | 88.31 |
| HMM + Maxent [48] | 88.76 |
| Semi-supervised [49] | 89.31 |
| Conv-CRF + Senna [44] | 89.59 |
| BLSTM-CRF [23] | 90.10 |
| CRF + LIE [50] | 90.90 |
| BLSTM-CRF (Ours) | **91.63** |

TABLE IV: Comparison of F1 scores of different models for NER

dependencies and achieve very good performance [39], [43], [45]. However, none of them surpass the performance of [46] which uses an HMM to capture the dependencies and a voting scheme to increase the confidence interval of the model. Our model outperforms all the existing models and achieves a state of the art F1 score of **96.76%**.

Table IV shows the results of the existing models on the NER task. Huang et al. [23] did many experiments using random and pre-trained embeddings on their model. For random embeddings, they achieved a very low score of 84.26%. However, when they use pre-trained SENNA embeddings [44] along with a gazetteer feature, their F1-score jumped up to 90.10% surpassing the Conv-CRF model [44] which uses window and sequence approach networks to do the tagging. Our model achieves a state of the art result of **91.63%**.

Table V shows the performance of our architecture in comparison with some top performing ones for the POS tagging task. As can be seen, a number of models use Convolution or LSTM or BLSTM to get the contextual features and CRF to do the tagging. They achieve very good accuracies of 97.29% [44], 97.51% [24] and 97.55% [27]. Some of the models use multitask learning, doing two or more tasks at the same time. They also achieve very good accuracies: 97.43% [28] and **97.59%** [25]. Our model achieves an accuracy of **97.58%** which is higher than all of the existing models except LM-LSTM-CRF [25] which leverages a language model for the tagging tasks. LM-LSTM-CRF, however, has a mean accuracy of **97.53%** (reported accuracy: $97.53 \pm 0.03$) which is lower than the our model's mean accuracy ($97.57 \pm 0.01$). Also, as

| Model | Accuracy |
|---|---|
| Conv-CRF [44] | 97.29 |
| 5wShapesDS [17] | 97.32 |
| Structure regularization [51] | 97.36 |
| Multitask learning [28] | 97.43 |
| Nearest neighbor [52] | 97.50 |
| LSTM-CRF [24] | 97.51 |
| LSTM-CNN-CRF [27] | 97.55 |
| LM-LSTM-CRF [25] | 97.59 |
| BLSTM-CRF (Ours) | 97.51 |
| BLSTM-CRF (Ours) without CNN | **97.58** |

TABLE V: Comparison of Accuracy of different models for POS tagging

| Model no. | Word emb | Sense | SP-CLSTM | Bigram | Suffix | | Spelling | | | Char Embed | | Prev. POS | Acc. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | CW | CO | R | CW | CO | CW | CO | | |
| 1 | Rand. | ✗ | ✗ | ✗ | - | - | - | - | - | - | - | ✗ | 95.42 |
| 2 | Glove | ✗ | ✗ | ✗ | - | - | - | - | - | - | - | ✗ | 96.13 |
| 3 | Glove | ✓ | ✗ | ✗ | - | - | - | ✓ | - | - | ✓ | ✗ | 97.08 |
| 4 | Glove | ✓ | ✓ | ✗ | ✓ | - | ✓ | - | - | ✓ | - | ✗ | 97.15 |
| 5 | Glove | ✓ | ✓ | ✗ | ✓ | - | - | ✓ | - | ✓ | - | ✗ | 97.22 |
| 6 | Glove | ✓ | ✗ | ✗ | - | ✓ | - | - | ✓ | ✓ | - | ✗ | 97.32 |
| 7 | Glove | ✓ | ✓ | ✗ | - | - | - | - | ✓ | ✓ | - | ✗ | 97.45 |
| 8 | Glove | ✓ | ✓ | ✗ | ✓ | - | - | - | ✓ | ✓ | - | ✗ | 97.48 |
| 9 | Glove | ✓ | ✓ | ✓ | ✓ | - | - | - | ✓ | ✓ | - | ✗ | 97.50 |
| 10 | Glove | ✓ | ✓ | ✗ | ✓ | - | - | - | ✓ | ✓ | - | ✓ | **97.58** |

TABLE VI: Ablation study of our BLSTM-CRF model for POS tagging. (R - Residual connection, CW - Concatenate with word embedding, CO - Concatenate with second last output layer)

| | Model | Acc. | Time |
|---|---|---|---|
| [25] | LSTM-CRF | 97.35 | 37 |
| | LSTM-CNN-CRF | 97.42 | 21 |
| | LM-LSTM-CRF | 97.53 | 16 |
| | LSTM-CRF | 97.44 | 8 |
| | LSTM-CNN-CRF | 96.98 | 7 |
| Ours | BLSTM-CRF | 97.58 | 4 |
| | BLSTM-CRF without CNN | 97.51 | 3.5 |

TABLE VII: Training time (in hours) of our BLSTM-CRF model on the WSJ corpus compared with all the models of [25] using the same hardware configuration (GPU: Nvidia GTX 1080)

| Module | 100th epoch | 200th epoch | 300th epoch | 400th epoch | 500th epoch |
|---|---|---|---|---|---|
| Word emb | 0.91 | 0.84 | 0.80 | 0.77 | 0.78 |
| Sense | 0.85 | 0.76 | 0.69 | 0.64 | 0.65 |
| SP-CLSTM | 0.81 | 0.66 | 0.48 | 0.35 | 0.34 |
| Bigram | 0.75 | 0.49 | 0.27 | 0.01 | 0.01 |

TABLE VIII: Change in $w$'s for each module with epochs.

shown in Table VII, our model's training time is one quarter that of LM-LSTM-CRF with on par performance.

Table VI gives the ablation study of our model where we show how we apply different combinations of features in different parts of our BLSTM-CRF architecture to get an optimal configuration. With so many features and parameters, these sequence models are very much prone to overfit. But with careful tuning as well as with proper feature connections, it is possible to leverage those features. We extract a set of morphological as well as semantic features from our dataset such as spelling, suffix and char-level features. We experiment on applying various combinations of these features in different segments of our model. Our extensive experimentation shows that optimal results are achieved when these features are added in the model through residual connection, concatenation with word embeddings and concatenation with the second last output layer. Focusing on which segment to connect each feature, our experiments found that the spelling feature works best when concatenated with the second last output layer,

and the suffix feature as well as the character embeddings work well when concatenated with the word embeddings. This configuration is what is kept in our final model. We further continue our experiments by turning on / off different modules such as word embedding, sense embedding, selective pickup from LSTM and bi-gram embedding. We found that the contribution of word embeddings, sense embeddings and selective pickup from LSTM are significant compared to the bigram module as shown by the weights at the 500th epoch in Table VIII. The bigram module works better whenever we do not consider the previously generated part-of-speech. So we kept the first three modules and discarded the bigram module from our final model. Our best model as shown in the last row of Table VI gives state of the art results.

## VI. CONCLUSION

In this paper, we propose an improved neural sequence labelling architecture by leveraging from additional linguistic information such as polysemy, bigrams, character level knowledge and morphological features. Benefitting from such adequately captured linguistic information, we can assemble a considerably more compact model, hence yielding much better training time without loss of effectiveness. To avoid feature collision we performed an extensive ablation study where we produced an optimal model structure along with an optimal set of features. Our best model achieved state of the art results on the POS tagging, NER and chunking benchmark datasets and at the same time remains four times faster to train than the best performing model currently available. Our experimental results show that multiple linguistic features and their proper inclusion significantly boosted our model performance.

## REFERENCES

[1] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 160–167.

[2] J. Li, M.-T. Luong, D. Jurafsky, and E. Hovy, "When are tree structures necessary for deep learning of representations?" *arXiv preprint arXiv:1503.00185*, 2015.

[3] X. Zheng, H. Chen, and T. Xu, "Deep learning for chinese word segmentation and pos tagging," in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 2013, pp. 647–657.

[4] L. Ratinov and D. Roth, "Design challenges and misconceptions in named entity recognition," in *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*. Association for Computational Linguistics, 2009, pp. 147–155.

[5] R. K. Ando and T. Zhang, "A high-performance semi-supervised learning method for text chunking," in *Proceedings of the 43rd annual meeting on association for computational linguistics*. Association for Computational Linguistics, 2005, pp. 1–9.

[6] D. Yarowsky, "Decision lists for lexical ambiguity resolution: Application to accent restoration in spanish and french," in *Proceedings of the 32nd annual meeting on Association for Computational Linguistics*. Association for Computational Linguistics, 1994, pp. 88–95.

[7] A. Ekbal, S. Mondal, and S. Bandyopadhyay, "Pos tagging using hmm and rule-based chunking," *The Proceedings of SPSAL*, vol. 8, no. 1, pp. 25–28, 2007.

[8] G. Zhou and J. Su, "Named entity recognition using an hmm-based chunk tagger," in *proceedings of the 40th Annual Meeting on Association for Computational Linguistics*. Association for Computational Linguistics, 2002, pp. 473–480.

[9] J. Lafferty, A. McCallum, and F. C. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," 2001.

[10] N. Nguyen and Y. Guo, "Comparisons of sequence labeling algorithms and extensions," in *Proceedings of the 24th international conference on Machine learning*. ACM, 2007, pp. 681–688.

[11] J. A. Benediktsson, P. H. Swain, and O. K. Ersoy, "Neural network approaches versus statistical methods in classification of multisource remote sensing data," 1990.

[12] N. Kalchbrenner, E. Grefenstette, and P. Blunsom, "A convolutional neural network for modelling sentences," *arXiv preprint arXiv:1404.2188*, 2014.

[13] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[14] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[15] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd international conference on Machine learning*. ACM, 2006, pp. 369–376.

[16] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in neural information processing systems*, 2014, pp. 3104–3112.

[17] C. D. Manning, "Part-of-speech tagging from 97% to 100%: is it time for some linguistics?" in *International Conference on Intelligent Text Processing and Computational Linguistics*, 2011, pp. 171–189.

[18] R. Sennrich and B. Haddow, "Linguistic input features improve neural machine translation," *arXiv preprint arXiv:1606.02892*, 2016.

[19] G. A. Miller, "Wordnet: a lexical database for english," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.

[20] P. Basile, A. Caputo, and G. Semeraro, "An enhanced lesk word sense disambiguation algorithm through a distributional semantic model," in *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, 2014, pp. 1591–1600.

[21] S. Banerjee and T. Pedersen, "An adapted lesk algorithm for word sense disambiguation using wordnet," in *International Conference on Intelligent Text Processing and Computational Linguistics*. Springer, 2002, pp. 136–145.

[22] S. Arora, Y. Li, Y. Liang, T. Ma, and A. Risteski, "Linear algebraic structure of word senses, with applications to polysemy," *arXiv preprint arXiv:1601.03764*, 2016.

[23] Z. Huang, W. Xu, and K. Yu, "Bidirectional lstm-crf models for sequence tagging," *arXiv preprint arXiv:1508.01991*, 2015.

[24] G. Lample, M. Ballesteros, S. Subramanian, K. Kawakami, and C. Dyer, "Neural architectures for named entity recognition," *arXiv preprint arXiv:1603.01360*, 2016.

[25] L. Liu, J. Shang, F. Xu, X. Ren, H. Gui, J. Peng, and J. Han, "Empower sequence labeling with task-aware neural language model," *arXiv preprint arXiv:1709.04109*, 2017.

[26] X. Yu, A. Faleńska, and N. T. Vu, "A general-purpose tagger with convolutional neural networks," *arXiv preprint arXiv:1706.01723*, 2017.

[27] X. Ma and E. Hovy, "End-to-end sequence labeling via Bi-directional LSTM-CNNs-CRF," *arXiv preprint arXiv:1603.01354*, 2016.

[28] M. Rei, "Semi-supervised multitask learning for sequence labeling," *arXiv preprint arXiv:1704.07156*, 2017.

[29] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional lstm and other neural network architectures," *Neural Networks*, vol. 18, no. 5-6, pp. 602–610, 2005.

[30] T. Thireou and M. Reczko, "Bidirectional long short-term memory networks for predicting the subcellular localization of eukaryotic proteins," *IEEE/ACM transactions on computational biology and bioinformatics*, vol. 4, no. 3, 2007.

[31] C. Dyer, M. Ballesteros, W. Ling, A. Matthews, and N. A. Smith, "Transition-based dependency parsing with stack long short-term memory," *arXiv preprint arXiv:1505.08075*, 2015.

[32] S. Bartunov, D. Kondrashkin, A. Osokin, and D. Vetrov, "Breaking sticks and ambiguities with adaptive skip-gram," in *Artificial Intelligence and Statistics*, 2016, pp. 130–138.

[33] T. S. Ferguson, "A bayesian analysis of some nonparametric problems," *The annals of statistics*, pp. 209–230, 1973.

[34] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.

[35] T. Lei, R. Barzilay, and T. Jaakkola, "Molding cnns for text: non-linear, non-consecutive convolutions," *arXiv preprint arXiv:1508.04112*, 2015.

[36] A. McCallum, D. Freitag, and F. C. Pereira, "Maximum entropy markov models for information extraction and segmentation." in *Icml*, vol. 17, no. 2000, 2000, pp. 591–598.

[37] A. Ratnaparkhi, "A maximum entropy model for part-of-speech tagging," in *Conference on Empirical Methods in Natural Language Processing*, 1996.

[38] S. Johansen and K. Juselius, "Maximum likelihood estimation and inference on cointegrationwith applications to the demand for money," *Oxford Bulletin of Economics and statistics*, vol. 52, no. 2, pp. 169–210, 1990.

[39] F. Sha and F. Pereira, "Shallow parsing with conditional random fields," in *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*. Association for Computational Linguistics, 2003, pp. 134–141.

[40] S. Bartunov, "Adagram," https://github.com/sbos/AdaGram.jl, 2017.

[41] T. Kudoh and Y. Matsumoto, "Use of support vector learning for chunk identification," in *Proceedings of the 2nd workshop on Learning language in logic and the 4th conference on Computational natural language learning-Volume 7*. Association for Computational Linguistics, 2000, pp. 142–144.

[42] T. Kudo and Y. Matsumoto, "Chunking with support vector machines," in *Proceedings of the second meeting of the North American Chapter of the Association for Computational Linguistics on Language technologies*. Association for Computational Linguistics, 2001, pp. 1–8.

[43] R. McDonald, K. Crammer, and F. Pereira, "Flexible text segmentation with structured multilabel classification," in *Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, 2005, pp. 987–994.

[44] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, "Natural language processing (almost) from scratch," *Journal of Machine Learning Research*, vol. 12, no. Aug, pp. 2493–2537, 2011.

[45] X. Sun, L.-P. Morency, D. Okanohara, and J. Tsujii, "Modeling latent-dynamic in shallow parsing: a latent conditional model with improved inference," in *Proceedings of the 22nd International Conference on Computational Linguistics-Volume 1*. Association for Computational Linguistics, 2008, pp. 841–848.

[46] H. Shen and A. Sarkar, "Voting between multiple data representations for text chunking," in *Conference of the Canadian Society for Computational Studies of Intelligence*. Springer, 2005, pp. 389–400.

[47] H. L. Chieu and H. T. Ng, "Named entity recognition: a maximum entropy approach using global information," in *Proceedings of the 19th international conference on Computational linguistics-Volume 1*. Association for Computational Linguistics, 2002, pp. 1–7.

[48] R. Florian, A. Ittycheriah, H. Jing, and T. Zhang, "Named entity recognition through classifier combination," in *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*. Association for Computational Linguistics, 2003, pp. 168–171.

[49] R. K. Ando and T. Zhang, "A framework for learning predictive structures from multiple tasks and unlabeled data," *Journal of Machine Learning Research*, vol. 6, no. Nov, pp. 1817–1853, 2005.

[50] A. Passos, V. Kumar, and A. McCallum, "Lexicon infused phrase embeddings for named entity resolution," *arXiv preprint arXiv:1404.5367*, 2014.

[51] X. Sun, "Structure regularization for structured prediction," in *Advances in Neural Information Processing Systems*, 2014, pp. 2402–2410.

[52] A. Søgaard, "Semisupervised condensed nearest neighbor for part-of-speech tagging," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, 2011, pp. 48–52.