

# Decoding a Generalized Neural Architecture for assigning labels to words

(This is a school project and not for ACL submission!!!)

## Anonymous ACL submission

### Abstract

This project is about a generalized neural network architecture that could be used for assigning labels to words in NLP tasks. The proposed system will be able to assign NER labels, POS Tags and can chunk sentences given a small training set. This implementation is based on the paper by May and Hovy ,2016.

### 1 Introduction

Most of the state of art systems currently in practice are specifically tailored for a single task with hand engineered features or neural networks specific for one task. Collobert et al (2011) called for and proposed a centralized architecture. Having a centralized architecture saves on development and debugging cost. There have been several improvements on this architecture.

### 2 Goals

The Goals of this project are two fold :

1. Decrease the run time of the implementation the authors have in theano / lasagne. The library of choice is tensorflow
2. provide insights on what each layer in the model does thereby increasing the interpretability of the model

This current project is a slight variation of one such papers (Ma & Hovy , 2016) which came up with a LSTM-CNN-CRF based system. The authors used a Lasagne based code and one the contributions of this paper is a tensorflow based implementation of the similar system. Preliminary studies of both implementation on a same system shows the tensor-flow implementation to be 2 times faster than the published theano version.

Another contribution of this project is the explanation it provides in understanding the components of the model. The model will be explained in terms of

1. why it assigned a label to the word in simple terms
2. understanding the role of CNN used for Character embedding
3. Understand the role of the word embedding layer

### 3 Prior work

There has been a couple of papers in the recent past using varying combinations of CNN or RNN with LSTM and CRF. One of the earliest papers in this area is from Baidu (Zhou & Xu. 2015) which is builds on Collobert et al.,2011 where they discuss a NN architecture with word em-bedding, CNN and CRF layer. The word em-beddings address the data sparsity problem while CRF is commonly used method in su-pervised slot tagging applications. However, long range dependencies are not modelled sufficiently by CNN which only includes words within a limited context. Some research groups (Melamud et. al, 2015 and Levy & Goldberg,2015) attempted to solve long range dependency problem by using a dependency parser. The authors from Baidu propose using a deep BiLSTM (i.e. 2 layers of BiLSTM) followed by a CRF. LSTM was selected because it can efficiently handle vanishing/exploding gradient while allowing words further apart in the sentence to influence the current word. The research groups which used dependency parser later published an-other paper where they created context2vec (Melamud et. al., 2016). context2vec has a BiLSTM followed by a MLP (multi-layer perceptron). Context2Vec helps

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099

100 identify words that fit in a similar context. Another interesting work is a paper from CMU (Ma  
 101 I& Hovy, 2016). The authors de-scribe an end-to-  
 102 end system without need for handcrafted feature  
 103 in detail.

## 104 105 106 107 4 Task Description & Review of existing literature

108 Here is a high-level summary of the tasks involved  
 109 with a focus on named entity recognition. While  
 110 considering entity tagging problems, here are a  
 111 couple of intuitions that will help us in building  
 112 a good language model.

- 113 114 115 116 117 118 1. Morphological information: (like the prefix  
 or suffix of word) provides pre-liminary hint  
 on the type of entity a word is. Ma and Hovy  
 (2016) used CNN over characters of a word  
 to model this.
- 119 120 121 122 123 124 125 126 127 128 129 130 131 2. Names often span across words so modeling  
 must be done jointly across multiple tokens (eg. Substrate Intelligence). Linear  
 CRF is often a method of choice (Ma and  
 Hovy,2016; Lample et al., 2016; Zhou &  
 Xu, 2015). It uses the transition probability  
 (i.e. probability of I-ORG given B-ORG)  
 and emission probability (i.e probability of  
 the word given the tag I-ORG) to generate  
 output tags with maximum likelihood. By us-  
 ing LC-CRF, NER grammar rules like I-ORG  
 cannot follow I-PER or I-PER cannot follow  
 B-LOC are implicitly encoded in the model
- 132 133 134 135 136 137 138 139 140 141 142 143 3. Orthographic evidence: This refers to the  
 internal characteristic of something that looks  
 like a name. For example, in some of  
 our project names (most likely abbreviations)  
 like QAS or LUNA we have u following L  
 and a following Q. In non-personal words, Q  
 is almost often followed by u. Lample et al  
 (2016) modeled this by using Bi-LSTM at the  
 character level. This will help us in modeling  
 words which are outside our known vocabu-  
 lary even when we have very little context for  
 them.
- 144 145 146 147 148 149 4. Distributional evidence: As Firth men-  
 tioned way back in 1950s, a word is known by  
 the context it keeps. Over the last year,  
 most p-pers (Melamud et al., 2016; Ma &  
 Hovy,2016; Lample et al., 2016; Zhou &

Xu, 2015) have used Bi-LSTM almost unani-  
 mously to address this issue. BiLSTM pro-  
 vides past and future context to a word.

- 150 5. Data sparsity: is another issue which has  
 151 been addressed using pretrained embedding  
 152 and by training new embeddings specific for  
 153 the data in hand.  
 154 6. Other features like Capitalization can also be  
 155 a good indicator of an entity.

156 However, not all these characteristics are always  
 157 present in every sample. Some names can sound  
 158 like regular words eg., keyphrase Identifica-  
 159 tion. Some of our project names are fully in capitals  
 160 while others may have just the first letter capital-  
 161 ized. So, many papers (Ma & Hovy,2016; Lample  
 162 et al., 2016; Zhou & Xu, 2015) report using drop  
 163 out layers to address this issue and to solve the  
 164 problem of overfitting.

## 165 166 167 168 5 Methods

### 169 170 171 5.1 Networks Architecture

172 The high level system architecture can be summa-  
 173 rized by the diagram from Ma I& Hovy, 2016.The  
 174 method used is similar to the above paper with  
 175 some minor variations. Dataset: ConLL 2003  
 176 NER dataset with 14041 / 3247 / 3450 sentences  
 177 in train / dev / test.

178 Word Embedding: We read in 6 billion pre-  
 179 trained vectors from Wikipedia and web (open  
 180 source).For every word in the vocabulary from  
 181 training set, if it does not have a pretrained vec-  
 182 tor randomly initialize it.Word dimension size is  
 183 set to 100

184 Character representations: We generate random  
 185 embeddings for every character in the vocabu-  
 186 lary. Embedding dimensions are set to 30. We  
 187 add a drop out layer with  $p=0.5$  to prevent over  
 188 fitting.Character representations need to be con-  
 189 verted to word format so we pass them into a  
 190 Convolutional Neural Network (CNN) with a filter  
 191 length =30 and window size of 3.Activation func-  
 192 tion is tanh.We use conv 1 D , activation as tanh,  
 full padding

193 Merge Layer: Word Embeddings and Character  
 194 representations are merged together and then se-  
 195 lected with a dropout rate of 0.5

196 BiLSTM:Use the standard LSTM architecture  
 197 as de-scribed by Hchreiter & Schmidhuber , 1997.  
 198 LSTM hidden layer dimension are set at 200.A  
 199 third drop out layer with  $p =0.5$  is added

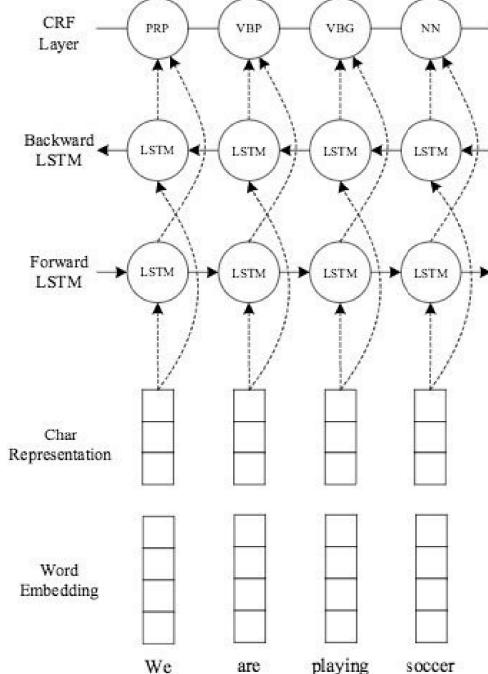


Figure 1: Model Architecture taken from Ma & Hovy, 2016.

## 5.2 Network shape

The network is composed of 3 layers. Each layer in terms of its output shape is described below(The input and out shape for each layer is described int he code):

- Layer 1 : This layer consists of character representations and word embeddings.
  - \* Character Representation layer: [batch\_size, max\_char\_per\_word \* sequence\_length, char\_embedd\_dim ]
  - \* Word Embedding layer :[batch\_size, sequence\_length, word\_embedd\_dim]
- Layer 2 : Forward and Reverse LSTMs (BiLSTM) along with an output linear dense layer to squash the results
  - \* Right LSTM : [batch\_size, sequence\_length, Hidden\_Unit\_Size]
  - \* Left LSTM : [batch\_size, sequence\_length, Hidden\_Unit\_Size]
  - \* BiLSTM : [batch\_size, sequence\_length, 2 \* Hidden\_Unit\_Size]
  - \* Linear Layer : [batch\_size, sequence\_length]

- Layer 3 : CRF Layer 250
- \* CRF : [batch\_size, sequence\_length] 251

## 5.3 Experimental Variables

Following experimental variables were changed during the course of experimentation:

- Gradient Clipping : Initially, there was an exploding gradient issue with the network returning NaNs after few runs. I tried Global gradient clipping ( absolute Vs Clip by Norm) , clipping after BiLSTM. I found absolute clipping after BiLSTM to give the best results. So it was used in subsequent experiments. 257  
258  
259  
260  
261  
262  
263  
264
- Dropout: The paper talked about 3 dropout layers interspersed within the network. I tried adding an additional dropout layer after BiLSTM and the accuracy dropped by more than 10 points. 265  
266  
267  
268  
269  
270
- Optimizer: The paper used SGD and mentioned that other optimizers gave similar results. I struggled with SGD for a long time and then just changed it Adam Optimizer. That small changes greatly increased the accuracy 271  
272  
273  
274  
275  
276
- Batch Size: Good results were obtained with small batch size i.e between 10 and 64 277  
278
- Filter Size and Number of Filters in CNN for character representation: Changing these values did not change the P/R numbers 280  
281  
282

## 5.4 Interpretability

The context which contributed to assigning particular label to a given word was obtained by modifying

- Ribeiro et. al., 2016 code to accommodate CRF function. 288  
289
- Adapting the CRF Viterbi decode to output probability (like) scores. 290  
291

The algorithm takes a sentence and creates variations by randomly removing words as shown in the Figure 2. Each of the variations is passed through the prediction function. The prediction function reads in the sentence and the target word and outputs the probability that the target word belongs to one of the 18 classes. We have 17 classes in the

- 300 • Sentence: Australian Tom Moody took six for 82 but Chris Adams , 123 and Tim  
 301 O'Gorman , 109 , took Derbyshire  
 302 • Variations(5000):  
 303   • Australian Tom Moody took six for 82 but Chris Adams , 123 and Tim O'Gorman , 109 ,  
 304    took Derbyshire  
 305   • Australian Tom Moody took six for 82 but Chris Adams , 123 and Tim O'Gorman , 109 ,  
 306    took Derbyshire  
 307   • Australian Tom Moody took six for 82 but Chris Adams , 123 and Tim O'Gorman , 109 ,  
 308    took Derbyshire

Figure 2: Variations of a sentence for predictions

309 training set and we add an extra unknown class  
 310 when we created the model. Now, output from  
 311 the predictor and the sentence is passed with the  
 312 output as label to the ridge regressor. The ridge  
 313 regressors assigns weights to each of the context  
 314 words

## 6 Evaluation

315 Extrinsic evaluation we are interested in is the F1  
 316 score. The F1 score is calculated as described by  
 317 Sang & Meulder , 2003. A named Entity is correct  
 318 if it is an exact match of the true entity.

$$F_{\beta} = \frac{(\beta^2 + 1) * precision * recall}{\beta^2 * precision + recall} \quad (1)$$

322 The intrinsic loss that the neural network optimizes  
 323 is the negative log likelihood of the CRF  
 324 loss as described by Ma & Hovy,2016.

## 7 Results and Discussion

### 7.1 Tensorflow Implementation

325 The Tensorflow implementation of the architecture  
 326 takes about 8 hours as opposed to 36 hours  
 327 for the theano implementation on the same  
 328 machine. I tried using both Gradient descent and  
 329 Adam Optimizers and found the Adam Optimizer  
 330 works better than stochastic gradient descent.  
 331 Also,I clipped the LSTM layer instead  
 332 of the whole network. I tried changing other  
 333 parameters like the learning rate and decay rate  
 334 and found the values suggested in the paper to  
 335 be most optimal. The whole implementation  
 336 can be obtained from <https://github.com/>  
 337 LopezGG/NN\_NER\_tensorFlow

#### 7.1.1 Precision & Recall

338 The overall statistics on the test set with break  
 339 down per category is

	ORG	LOC	MISC	PER	Other
Precision	57.72	74.41	62.69	69.82	97.19
Recall	52.18	77.39	63.89	66.94	98.04
F1 score	54.81	75.87	63.28	68.34	97.61

## 7.2 Interpretability

### 7.2.1 Character Embedding layer

350 Here are some of the words which were not in the  
 351 training and development set vocabulary. These  
 352 words are called Out of Unsupervised Vocabulary

9-0-54-1	disapproval	highest-ranking
9-0-49-1	disgraceful	bigest-ever
6-0-30-1	scrapped	Rochester
6-0-40-1	distracted	rapidly-growing
4-0-18-0	dishoarding	higher-than-expected
8-0-37-2	non-partisan	self-government

353 The character embedding layer groups words  
 354 which share a similar pattern as shown above. This  
 355 helps with categorizing unknown words

### 7.2.2 Word Embedding layer

356 The original pretrained word embeddings were  
 357 grouped together by their semantic meaning. Af-  
 358 ter running those pretrained words through this  
 359 network , the Word embedding layer now is op-  
 360 timized to group words belonging to the same cat-  
 361 egory as shown in the table below

Commission:E-ORG	Democratic:S-MISC
Assembly:E-ORG	GMT:S-MISC
Dortmund:E-ORG	English:S-MISC
University:E-ORG	Moslem:S-MISC
Ministry:E-ORG	Thai:S-MISC
TVM:S-ORG	Polish:S-MISC

362 It is to be noted that many of the words belong to  
 363 multiple categories. for eg. the first word "Com-  
 364 mission: E-ORG: 32; S-ORG: 11; B-MISC: 4; I-  
 365 ORG: 2 " belongs to E-ORG category in 32 in-  
 366 stances. The categories to which the word "Com-  
 367 mission" belongs to depends on the context. The  
 368 table above shows the most frequent class assigned  
 369 to each word.

### 7.2.3 Debugging the CRF layer

370 The results for interpreting the CRF layers are as  
 371 follows:

- 372 • Sentence: Australian Tom Moody took six  
 373 for 82 but Chris Adams , 123 and Tim  
 374

400	O’Gorman , 109 , took Derbyshire The target word is ”O’Gorman” with a tag E-PER which the system has never seen before in the development or training set. The model shows that the presence of ”Tim” , B-PER had a high influence on assigning E-PER to ”O’Gorman”	450
401		451
402		452
403		453
404		454
405		455
406		456
407		457
408		458
409		459
410		460
411		461
412		462
413		463
414		464
415		465
416		466
417		467
418		468
419		469
420		470
421		471
422		472
423	The target word in this case was ”Indian” (E-MISC). the presence of ”West” (B-MISC) contributed to this assignment.	473
424		474
425		475
426		476
427	The actual scores from the models for all these sentences are shown in figures	477
428		478
429		479
430		480
431		481
432		482
433	<b>8 Challenges &amp; Opportunities</b>	483
434	We need a better way to visualize the results of LSTM and CRF. The current implementation is hacky because the CRF-decode does not actually output probability. I will have to read and explore other papers to make this calculation more robust	484
435		485
436		486
437		487
438	<b>9 References</b>	488
439	Felix A Gers, Jurgen Schmidhuber, and Fred Cummins. 2000. Learning to forget: Continual prediction with lstm. Neural computation, 12(10):24512471.	489
440		490
441		491
442		492
443	Jie Zhou and Wei Xu. 2015. End-to end learning of semantic role labeling using recurrent neural networks. In Proceedings of the Annual Meeting of the Association for Computational Linguistics.	493
444		494
445		495
446		496
447	Ronan Collobert, Jason Weston, Leon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. Journal of Machine Learning Research, 12:24932537, November	497
448	Tjong Kim Sang, Erik. F. 2002. Introduction to the CoNLL-2003 Shared Task: Language Inde-	498
449		499