

Speed versus Accuracy in Neural Sequence Tagging for Natural Language Processing

by

Xinxin Kou

B.Sc. (Hons.), Dalhousie University, 2015

Thesis Submitted in Partial Fulfillment of the
Requirements for the Degree of
Master of Science

in the
School of Computing Science
Faculty of Applied Sciences

© Xinxin Kou 2017
SIMON FRASER UNIVERSITY
Fall 2017

All rights reserved.

However, in accordance with the *Copyright Act of Canada*, this work may be reproduced without authorization under the conditions for “Fair Dealing.” Therefore, limited reproduction of this work for the purposes of private study, research, education, satire, parody, criticism, review and news reporting is likely to be in accordance with the law, particularly if cited appropriately.

Abstract

Sequence Tagging, including part of speech tagging and named entity recognition, is an important task in NLP. Recurrent neural network models such as Bidirectional LSTMs have produced impressive results on sequence tagging. In this work, we first present a simple and fast greedy sequence tagging system using different types of feedforward neural network models. Then we show the speed and accuracy comparison between Bidirectional LSTMs and feedforward models. Besides the feedforward and the Bidirectional LSTM models, we propose two new models based on Mention2Vec by Stratos (2016): Feedforward-Mention2Vec for Named Entity Recognition and BPE-Mention2Vec for Part-of-Speech Tagging. Feedforward-Mention2Vec predicts named entity boundaries first and then predicts types of named entities. BPE-Mention2Vec uses the Byte Pair Encoding algorithm to segment words in a sequence first and then predicts the Part-of-Speech tags for the subword spans. We carefully design the experiments to demonstrate the speed and accuracy trade-off in different models. The empirical results reveal that feedforward models can achieve comparable accuracy and faster speed than recurrent models for Part-of-Speech tagging, and Feedforward-Mention2Vec is competitive with the fully structured BiLSTM model for Named Entity Recognition while being more scalable in the number of named entity types.

Keywords: Natural Language Processing; Sequence Tagging, Deep Learning

Acknowledgements

I would like to express my profound sense of gratitude to my supervisor Dr. Anoop Sarkar for introducing me to this research topic and providing his continuous support and valuable guidance throughout my graduate study. I can not imagine having a better advisor and mentor. In addition, I would like to express my sincere appreciation to Dr. Fred Popowich for his useful advice and feedback on this work, and Dr. Jiannan Wang for being my defence examiner and reading my thesis.

Table of Contents

Abstract	ii
Acknowledgements	iii
Table of Contents	iv
List of Tables	vi
List of Figures	vii
1 Introduction	1
1.1 Sequence Tagging Task	1
1.1.1 Part-of-Speech Tagging (POS)	1
1.1.2 Named Entity Recognition (NER)	1
1.2 Motivation	3
1.3 Contribution	5
1.4 Overview	6
2 Feedforward Neural Network Models	7
2.1 Network Architecture	7
2.1.1 Feedforward-History Model	7
2.1.2 Feedforward-CRF	10
2.2 Experiments and Results	11
3 Bidirectional Long Short Term Memory Network Models	14
3.1 Model Description	14
3.1.1 Bidirectional Bidirectional Long Short Term Memory	14
3.1.2 Character Embeddings	16
3.2 Experiments and Results	16
4 Mention2Vec Models	19
4.1 Model Description	19
4.1.1 Feedforward-Mention2Vec for NER	19

4.1.2	BPE-Mention2Vec for POS	22
4.2	Experiments and Results	23
5	Discussion	28
6	Conclusion and Future Work	32
6.1	Contribution	32
6.2	Future Work	33
	Bibliography	34

List of Tables

Table 1.1	Number of Sentences (words) and Labels in Each Training, Validation and Test Section of Different Data Sets	2
Table 1.2	Performance of Different POS and NER Systems	2
Table 1.3	Name Entity Types in OntoNotes	3
Table 2.1	Hyperparameters used in Feedforward Models	11
Table 2.2	Feedforward Models Accuracy and F1 Score	12
Table 2.3	Feedforward Models Decoding Speed	12
Table 3.1	Hyperparameters used in BiLSTM Models	17
Table 3.2	BiLSTM Models Accuracy and F1 Score	18
Table 3.3	BiLSTM Models Decoding Speed	18
Table 4.1	F1 Score and Decoding Speed Comparison on CoNLL 2003	25
Table 4.2	F1 Score and Decoding Speed Comparison on OntoNotes	25
Table 4.3	Per-label F1 Score on OntoNotes	26
Table 4.4	POS Tagging System Accuracy and Decoding Speed	26
Table 5.1	Neural Network Models Accuracy and F1 Score	29
Table 5.2	Neural Network Models Decoding Speed	29
Table 5.3	F1 Scores and Decoding Speed on CoNLL 2003 and OntoNotes . . .	31

List of Figures

Figure 1.1	An example of POS	3
Figure 1.2	An example of NER	4
Figure 2.1	The architecture of a greedy tagging system using the Feedforward-History model.	8
Figure 2.2	Performance comparison between Feedforward Models on POS and NER (accuracy for POS; F1 score for NER	12
Figure 3.1	The architecture of an NER tagging system using BiLSTM with CRF	15
Figure 3.2	The word embedding derived from the character embeddings	17
Figure 3.3	Performance comparison between BiLSTM models on POS and NER (accuracy for POS; F1 score for NER	18
Figure 4.1	The first step of Mention2Vec for NER	20
Figure 4.2	The second step of Mention2Vec for NER	20
Figure 4.3	The third step of Mention2Vec for NER	21
Figure 4.4	The first step of Feedforward-Mention2Vec for NER	22
Figure 4.5	An example of using BPE for word segmentation	23
Figure 4.6	An example of using BPE-Mention2Vec to find POS tags	24
Figure 5.1	Results of the POS system using different Neural Network Models .	30
Figure 5.2	Results of the NER system using different Neural Network Models on CoNLL 2003	31

Chapter 1

Introduction

In this chapter, we first describe sequence tagging tasks and introduce the motivation of this thesis. Then, we summarize our major contributions and describe the structure of the thesis.

1.1 Sequence Tagging Task

1.1.1 Part-of-Speech Tagging (POS)

Part-of-Speech Tagging (POS) is a basic sequence tagging task, which assigns each word with a unique tag that indicates its syntactic role, such as noun, adverb, verb, etc. . . . Figure 1.1 illustrates a typical POS task.

Most POS systems are evaluated on the English Penn TreeBank data set (Marcus et al., 1993), which contains 45 part-of-speech tags. The standard split uses section 1-18 of the Treebank for training, section 19-21 for tuning, and section 22-24 for testing (Toutanova et al., 2003). The experimental data are summarized in Table 1.1. Many existing models are linear models, such as Maximum entropy Markov models (MEMMs) which obtain 96.46% per word accuracy (McCallum et al., 2000); and the averaged perception discriminative model which obtains 97.11% per word accuracy (Collins, 2002). More recently, neural network models have been proposed to improve the state-of-the-art score. The Bidirectional LSTM network model proposed by Huang et al. (2015) reaches 97.55% per word accuracy. Ling et al. (2015) presents the compositional character-to-word LSTM model which reaches the state-of-the-art performance: 97.78% per word accuracy. The performance of existing POS models is reported in Table 1.2.

1.1.2 Named Entity Recognition (NER)

Named Entity Recognition (NER) is a more complex sequence tagging task than POS, which identifies expressions that refer to named entities, such as peoples, places, organizations and

Table 1.1: Number of Sentences (words) and Labels in Each Training, Validation and Test Section of Different Data Sets

	Training	Validation	Test	labels
Penn Treebank	39831(950011)	1699(40068)	2415(56671)	45
CoNLL2003	14987(204567)	3466(51578)	3684(46666)	9
OntoNotes	75187(1088503)	9479(147724)	9603(152728)	18

Table 1.2: Performance of Different POS and NER Systems

POS Systems	Accuracy	NER Systems	F1
McCallum et al. (2000)	96.46%	Florian et al. (2003)	88.76
Collins (2002)	97.11%	Huang et al. (2015)	88.83
Huang et al. (2015)	97.55%	Lample et al. (2016) with CRF	90.94
Ling et al. (2015)	97.78%	Lample et al. (2016) with Stack LSTM	90.33

others. The main difference between NER and POS is that each named entity label can span multiple words while each part-of-speech tag is only for one word. A popular convention in NER is to use the "IOB" label scheme (Inside, Outside, Beginning): if the word is the beginning of a named entity label, it is marked with B-label; if the word is inside a named entity label but not the first one, it is marked with I-label; if the token is outside the named entity, it is marked with "O". An example of NER is shown in Figure 1.2. There are different numbers of named entity types in different NER tasks. For example, the shared task of CoNLL 2003 (Tjong Kim Sang and De Meulder, 2003) contains 4 types of named entities: locations (LOC), persons (PER), organizations (ORG), and miscellaneous (MISC); and the OntoNotes English data set (Hovy et al., 2006) contains 18 types of named entities shown in Table 1.3. The predefined training, development, and testing split of the data sets are shown in Table 1.1.

Most NER models are evaluated on CoNLL 2003 data set and are measured by the F1 score, which is the harmonic mean of precision and recall.

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (1.1)$$

Since CoNLL 2003 has a relatively smaller amount of training data, most of the existing models make use of pretrained word embeddings along with the training data. A commonly used pretrained word embedding is GloVe (Pennington et al., 2014) which contains 40K words. The best system presented at the NER CoNLL 2003 challenge by Florian et al. (2003) obtains 88.76 F1 score. The model using Bidirectional LSTM by Huang et al. (2015) reaches 88.83 F1 score. Both of these two models use many external features along with gazetteer¹. Lample et al. (2016) proposed two NER models with no external features or a gazetteer: the first one makes structured prediction using Bidirectional LSTM, Character

¹Gazetteers are language-specific knowledge resources

Table 1.3: Name Entity Types in OntoNotes

Types	Named Entities
PERSON	People, including fictional
NORP	Nationalities or religious or political groups
FACILITY	Buildings, airports, highways, bridges, etc.
ORG	Companies, agencies, institutions
GPE	Countries, cities, states
LOC	Non-GPE locations, mountain ranges, bodies of water
PRODUCT	Vehicles, weapons, foods, etc. (Not services)
EVENT	Named hurricanes, battles, wars, sports events
WORK OF ART	Titles of books, songs
LAW	Named documents made into laws
LANGUAGE	Any named language
DATE	Absolute or relative dates or period
TIME	Times smaller than a day
PERCENT	Percentage
MONEY	Monetary values, including unit
QUANTITY	Measurements, as of weight or distance
ORDINAL	First, Second, etc.
CARDINAL	Numerals that do not fall under another type

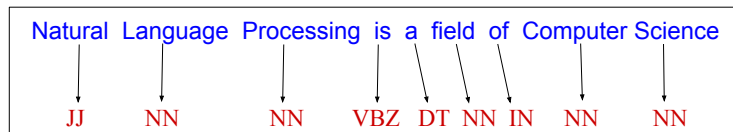


Figure 1.1: An example of POS

Embeddings (Ling et al., 2015) and Conditional Random Field (CRF) (Lafferty et al., 2001); and the second one uses a Shift-Reduce framework with Stack-LSTM (Dyer et al., 2015). The first model achieves the state-of-the-art F1 score, while the second one performs slightly worse. The performance of existing NER models is reported in Table 1.2.

1.2 Motivation

Recurrent Neural Networks (RNNs) have obtained impressive results in many NLP tasks, such as speech recognition (Graves et al., 2013) and machine translation (Cho et al., 2014). Bidirectional Long Short Term Memory (BiLSTM) (Hochreiter and Schmidhuber, 1997; Graves and Schmidhuber, 2005) is one of the RNN architectures that can maintain long-

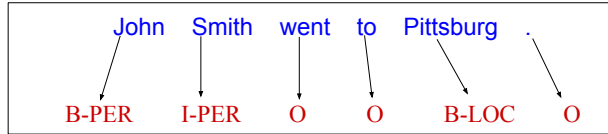


Figure 1.2: An example of NER

distance information from the past and future elements in a input sequence. Based on the existing work (POS system by Ling et al. (2015) and NER system by Lample et al. (2016)), state-of-the-art results on POS and NER can be obtained by using a BiLSTM network with Character Embeddings and CRF. Some work has also shown that feedforward neural networks can achieve comparable accuracy to recurrent models in tasks such as POS and Dependency Parsing (Andor et al., 2016). One approach presented in this thesis is to employ a greedy transition system with a feedforward network to make independent classification decisions on each word. However, the greedy system is limited when there are strong correlations between output labels. NER is one such task which has grammar constraints on output label sequences. For example, "I-PER" cannot follow "B-LOC" in NER. In order to take into account the strong dependencies among output labels, a CRF layer is added to model the output label sequence jointly. Since the CRF-based models focus on the sentence level and compute the score of every possible sequence, they take more time in training and decoding. Since we are interested in the decoding speed and performance trade-off in different neural network models, we re-implement variants of feedforward models and BiLSTM models using Tensorflow (Abadi et al., 2016) and systematically compare the performance and decoding speed between them on sequence tagging tasks, such as POS and NER.

Since the named entity labels in NER often span multiple tokens, most neural architectures for NER predict the boundaries and types of entities together using the "IOB" label scheme. Mention2Vec (Stratos, 2016) is proposed to address the natural segment-level representation in NER by separating the NER task into boundary detection (I, O, B) and type prediction (PER, LOC, etc.). While Mention2Vec employs two BiLSTMs for each sub-task, we replace the BiLSTM layer for boundary detection with a feedforward network in order to obtain a simpler model and accelerate the decoding process. This new model is denoted as Feedforward-Mention2Vec in this thesis.

We use Byte Pair Encoding (BPE) (Gage, 1994) to deal with rare words in sentences for machine translation (Sennrich et al., 2015), and we propose a new model combining BPE and Mention2Vec for POS, which is denoted as BPE-Mention2Vec in this thesis. We use BPE

to segment the input words in the hope of capturing the orthographic evidence of the words without using spelling features (like prefixes and suffixes) or Character Embeddings. After we segment input words, POS becomes an NER-like task. Then, we can use Mention2Vec to solve the rest of the problem. Since boundaries of output tags are known in POS, we only need to use a BiLSTM network for type prediction in BPE-Mention2Vec.

1.3 Contribution

The three main contributions of this thesis are:

1. We implement two greedy tagging systems with two different feedforward network models, Feedforward-History and Feedforward-CRF.

Feedforward-History takes word features in context, spelling features, and previous tag features as input. Feedforward-CRF uses word features in context and CRF to model the output sequences jointly. There is little existing work measuring the decoding time using feedforward networks on sequence tagging, so we conduct experiments on NER and POS and record the performance and decoding speed on Penn Treebank data for POS and CoNLL 2003 data for NER. To test the robustness of the feedforward networks, we also conduct experiments using a feedforward network with only word features, which also serves as the baseline. We compare different feedforward models and provide analysis on the results.

2. In addition to the greedy tagging systems with different feedforward networks, we build a tagging system using the fully structured BiLSTM model, which makes use of Character Embeddings and CRF. There is little work examining how the configurations of the fully structured BiLSTM model affect the decoding speed, such as whether the model is using Character Embeddings or the model is using CRF. We conduct experiments to compare the performance and decoding time of different configurations on POS and NER. We also benchmark the BiLSTM models against the feedforward models.
3. Last but not least, we introduce two new neural architectures based on Mention2Vec: Feedforward-Mention2Vec for NER and BPE-Mention2Vec for POS.

Originally, Mention2Vec was designed for NER and used BiLSTMs to detect named entity boundaries and predict corresponding types separately. We propose Feedforward-Mention2Vec for NER, in which we use a feedforward network with CRF to predict named entity boundaries instead of using a BiLSTM. We denote this new model for NER as Feedforward-Mention2Vec. We also adapt Mention2Vec for POS by combining Byte Pair Encoding (BPE) with BiLSTM in the model. BPE is used to segment the input words in our model, and it converts POS into a NER-like task. We denote

this new model for POS as BPE-Mention2Vec. In BPE-Mention2Vec, we use a feed-forward network to compute the hidden embeddings of the input segmented words. Since the boundaries of subword units are known, BPE-Mention2Vec does not need to predict the boundaries. It takes the hidden embeddings, tag boundaries, and a BiLSTM network to predict the actual POS tags. We benchmark these two multi-task models against the state-of-the-art BiLSTM model on POS and NER. Since different NER tasks have different numbers of named entity types, the decoding time of a fully structured BiLSTM model grows quadratically in the number of types. In Mention2Vec and Feedforward-Mention2Vec, we only apply CRF on boundary labels ('I', 'O', 'B'), so the decoding time grows linearly in the number of types. To show the time difference on different NER tasks, we conduct the NER experiments on CoNLL 2003 which contains 4 different named entity types and OntoNotes which contains 18 different named entity types.

1.4 Overview

The thesis is organized as follows:

In **Chapter 2** we present three different feedforward network models: feedforward network with only word features (denoted as Feedforward); feedforward network with spelling features and history features (denoted as Feedforward-History); and feedforward network with CRF (denoted as Feedforward-CRF). We explain the training and decoding process of the feedforward models, describe the experiment design, and compare the performance and decoding time between different feedforward models.

In **Chapter 3** we present the fully structured BiLSTM model (a combination of BiLSTM, Character Embeddings and CRF), denoted as the BiLSTM-Char-CRF model. We explain the training and decoding process of BiLSTM-Char-CRF, and conduct experiments using three different configurations: BiLSTM with word features only, BiLSTM with Character Embeddings, and BiLSTM model with Character Embeddings and CRF.

In **Chapter 4** we present two new new multi-task models based on Mention2Vec: Feedforward-Mention2Vec for NER and BPE-Mention2Vec for POS. We explain the architectures of these two models and we compare the performance and speed between multi-task models and other models in this thesis.

In **Chapter 5** we discuss the empirical results from the previous chapters. We also analyze the trade-off between performance and speed in different neural network models.

In **Chapter 6** we summarize the contributions of this thesis and discuss ongoing and related future work.

Chapter 2

Feedforward Neural Network Models

In this chapter, we describe two feedforward network models: Feedforward-History and Feedforward-CRF. Then, we show the performance and decoding speed of different feedforward models on POS and NER.

2.1 Network Architecture

2.1.1 Feedforward-History Model

Inspired by the greedy parser system by Chen and Manning (2014), we present a similar greedy transition system for sequence tagging in this thesis. The greedy parser system employs a basic arc-standard system (Nivre and Scholz, 2004), which consists of three types of transitions (LEFT-ARC, RIGHT-ARC, and SHIFT), a stack, and a buffer. While the greedy parser system has three types of actions, the sequence tagging system only needs the SHIFT action which predicts the tag of the current word in the buffer and shifts the word to the stack. Syntaxnet from Google also uses the same model for POS and dependency parsing and achieves near the state-of-the-art performance on both tasks (Alberti et al., 2017).

In our greedy tagging system, we use a feedforward network to make decisions on individual words, and we assume that the word to be labeled depends mainly on its neighbor words instead of the whole sentence. Besides the word features, the system also takes into account the lexical composition of the words (spelling features), and the previous tag decisions (previous tag features). Thereby, we represent this architecture as the Feedforward-History model. The way Feedforward-History incorporates word features is the same with the window approach proposed in Collobert et al. (2011). The difference between these two models is that Feedforward-History uses spelling features and previous tag features while the model in Collobert et al. (2011) only uses word features.

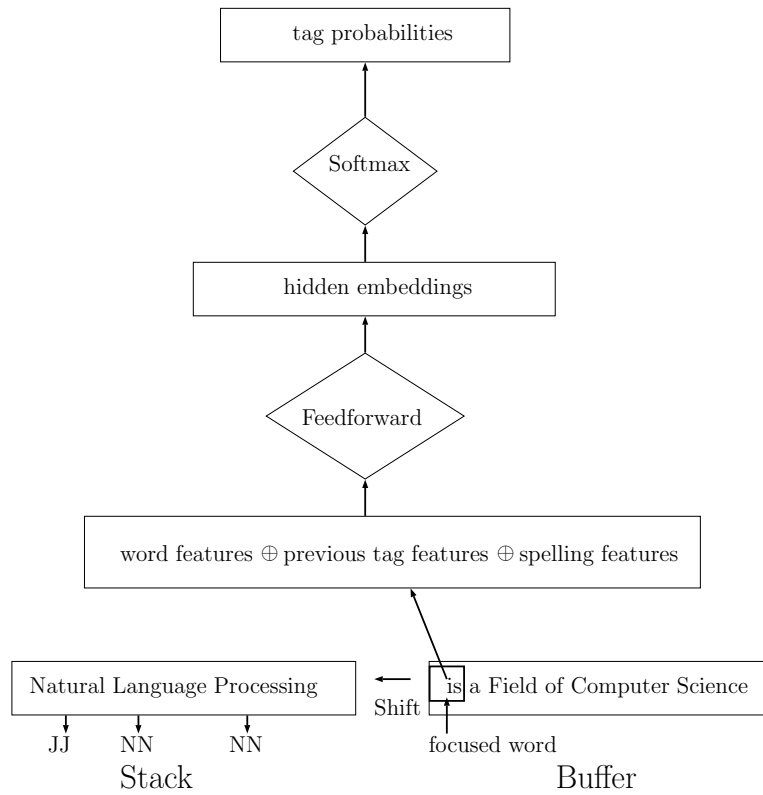


Figure 2.1: The architecture of a greedy tagging system using the Feedforward-History model.

Figure 2.1 illustrates the process of how a greedy tagging system uses Feedforward-History to decode a POS example. Since the focused word only depends on its neighbors in this greedy tagging system, we use a fixed size window around the current word to generate features. In order to generate features for the words at the beginning and at the end of the sentence, we add special padding words at the beginning and at the end. To generate dense word features, we convert each word in the input sequence to a d -dimensional vector representation e_{w_i} . Meanwhile, we have a full vocabulary embedding dictionary E_w . Given a word w_i , we look up its embedding e_{w_i} in E_w . According to Ratnaparkhi et al. (1996), spelling features of a word can help predict the part-of-speech tag. The examples of spelling features are upper and lower case features, prefix and suffix features, digit features, etc. Each spelling feature of a word is also associated with an embedding vector e_{s_i} , and it can be looked up from an embedding vector dictionary E_s . Besides word features and spelling features, we incorporate the previous tag features in this model. Each previous tag decision is represented as e_{ht_i} and can be looked up from a tag dictionary represented as E_{HT} .

The input layer $X = (x_1, x_2, \dots, x_n)$ to the feedforward network is obtained by concatenating the word feature vectors, spelling feature vectors, and previous tag feature vectors. In general, generating a lot of hand-engineered features for sequence tagging is expensive: selecting useful features is an empirical process based on trials and errors, and computing feature vectors requires searching feature strings in huge dictionaries and concatenating them together. We try to use as little hand-engineered features as possible to save the time from feature generation while keeping the model accurate. In our implementation, we extract the word and spelling features on a window size of 8 centered on the focused word. We extract the following spelling features of each word: whether it starts with a capital letter; whether it has all capital letters; whether it has a mix of letters and digits; whether it has punctuation; and letter prefixes and suffixes of length two and three. We also extract 4 previous tags as parts of input features.

While the input layer is the concatenation of the feature vectors of the focused word, the output layer is a probability distribution over tags. The input unit x_i is mapped to a hidden unit h_i through the rectifier activation function (*ReLU*) (Nair and Hinton, 2010):

$$ReLU(x) = \text{Max}(0, x) \quad (2.1)$$

$$h_i = ReLU(W_1^w x_i^w + W_1^s x_i^s + W_1^l x_i^l + b_1), \quad (2.2)$$

where x^w represents the word features, x^s represents the spelling features, x^l represents previous tag features, W_1 is the weight parameters in the hidden layer, and b_1 is the bias of in the hidden layer. The output of the network is a probability distribution over tags,

and its dimension is the size of all tags. The tag probability distribution is modeled by a softmax layer (*softmax*) (Dugas et al., 2001):

$$p_i = \text{softmax}(W_2 h_i + b_2), \quad (2.3)$$

where W_2 is the weight parameter in the softmax layer and b_2 is the bias in the softmax layer. The network is trained by minimizing a negative log likelihood over the training data. Embedding vectors are trained together with weight vectors and bias in the model. We denote all trainable parameters as θ . Given a sequence predictions, $Y = (y_1, y_2, \dots y_n)$, the score of the output sequence is the sum of the probability of each decision y_i :

$$S(Y) = \sum_i^n p_i(y_i), \quad (2.4)$$

and the loss function is:

$$L(\theta) = -\log \left(\sum_i^n p_i(y_i) \right), \quad (2.5)$$

2.1.2 Feedforward-CRF

There are two ways to incorporate output tags in the model: the first one is to use the previous predicted tags as input features to feed into the feedforward network as in Feedforward-History; the second one is to use CRF on sentence level output tags as in the following model. Feedforward-CRF shares the same architecture with Feedforward-History except using CRF to model the output sequence after computing tag probability distribution for each individual word. Feedforward-History can perform well on some sequence tagging tasks in which output labels do not have strong correlations, such as POS. Some tasks have grammar constraints on output tags so that there are strong dependencies among them, such as NER. While Feedforward-History fails to take into account the grammar constraints on the output tag sequences, Feedforward-CRF can make final decisions on the sentence level. Since the first part of the architecture is the same as Feedforward-History, we skip this part. We describe how to apply CRF to model the sequence in detail here.

Given a sequence of output predictions $y = (y_1, y_2, \dots y_n)$, the score of the output sequence is:

$$S(y) = \sum_i^n T_{i,i+1} + \sum_i^n p_i(y_i) \quad (2.6)$$

where T is a matrix of transition scores such that $T_{i,j}$ represents the score of a transition from the label i to label j , and $p_i(y_i)$ is the probability of y_i being the label of word i .

During training, we score every possible output sequence, and use a softmax layer to generate the probability distribution of output sequences, shown in Equation 2.7. Then, we

Table 2.1: Hyperparameters used in Feedforward Models

Hyperparameters	Values
window size	8
word embedding size	100
feedforward layer	1
feedforward layer size	200
optimizer	Adam
learning rate	0.01
batch size	32

minimize the negative log likelihood over the training sentences. During decoding, we can recover the tag sequences of test data using dynamic programming.

$$P(y) = \text{softmax}(S(y)) \quad (2.7)$$

2.2 Experiments and Results

In POS experiments, we train the models using the Penn Treebank training data and development data. Then, we decode the test data using trained models and record the per word accuracy and decoding time. In NER experiments, we train the models using the CoNLL 2003 training data and development data. Then we decode the test data using the trained models, and record the F1 score and decoding time. The details of the data sets are shown in Table 1.1. The performance of POS is measured by per word accuracy, and the performance of NER is measured by F1 score. The speed is measured by the average number of sentences and words decoded per second. We conduct experiments using Feedforward-History and Feedforward-CRF. We also want to show robustness of the feedforward network models, so we build a model using a feedforward network with only word features, which serves as the baseline model in our experiments. The tagging system using a feedforward network with only word features is denoted as Feedforward model in this thesis.

We implement the models using Python and the Tensorflow 1.0 library. In NER experiments, because of the small amount of training data, we use GloVE pretrained word embedding where each word corresponds to a 100-dimensional embedding vector. We use the hyper-parameters defined in the work of Chen and Manning (2014). Specifically we use Adam (Kingma and Ba, 2014) for stochastic optimization, and set the learning rate to be 0.01. We use 1 hidden layer and set the hidden layer size to be 200. We have a batch implementation which can process multiple sentences at the same time, and we set the batch size to be 32 in the experiments. Table 2.1 shows the hyperparameters. We initialize all the weight parameters using Xavier initialization (Glorot et al., 2011). To measure the speed, we run all the experiments in this thesis on a GeForce GTX 1080 GPU.

Table 2.2: Feedforward Models Accuracy and F1 Score

Model	POS (Accuracy)	NER (F1 Score)
Feedforward	95.89	84.12
Feedforward-History	97.28	86.54
Feedforward-CRF	97.30	87.85

Table 2.3: Feedforward Models Decoding Speed

Model	POS (sentences, words/sec)	NER (sentences, words/sec)
Feedforward	1319(30967)	2117(26819)
Feedforward-History	829(19474)	1390(17609)
Feedforward-CRF	761(17877)	1374(17412)

Comparison of Feedforward Models

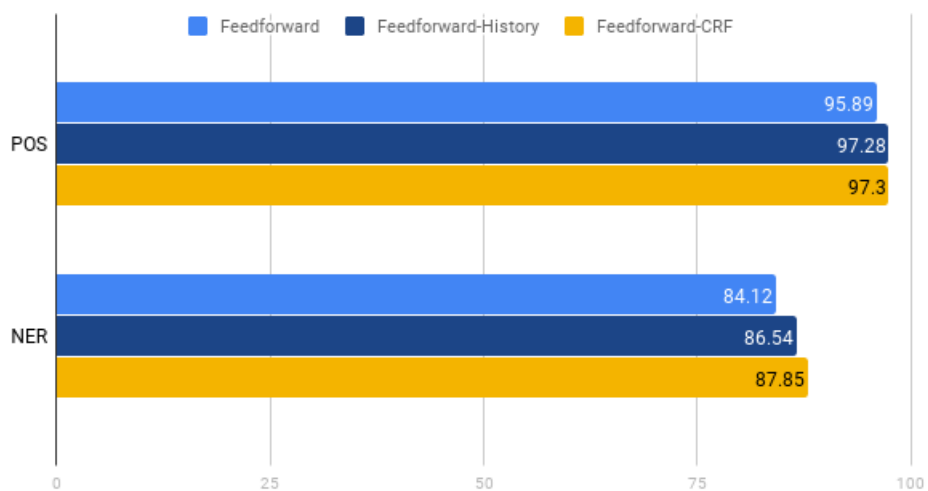


Figure 2.2: Performance comparison between Feedforward Models on POS and NER (accuracy for POS; F1 score for NER)

Table 2.2 and Table 2.3 present our final benchmark results of three feedforward models on POS and NER. Figure 2.2 shows the difference in the performance of the three models. In POS, Feedforward is 1.39% less accurate than Feedforward-History. In NER, the F1 score of Feedforward is 2.42 lower than Feedforward-History. Since Feedforward model with only word features does not decrease the performance dramatically, we can conclude that our greedy tagging systems with feedforward models do not heavily rely on spelling features. It also reveals that the spelling features are more helpful in NER than in POS. In both POS and NER, Feedforward-CRF has better performance than Feedforward-History, but it improves the NER performance more because of the dependencies among output tags. Feedforward-CRF is slower than Feedforward-History since CRF computes the scores of every possible output sequences.

Chapter 3

Bidirectional Long Short Term Memory Network Models

In this chapter, we first describe Bidirectional Long Short Term Memory network (BiLSTM) and Character Embeddings. Then, we show the performance and decoding speed of BiLSTM models with different configurations on POS and NER.

3.1 Model Description

3.1.1 Bidirectional Bidirectional Long Short Term Memory

One major disadvantage of feedforward neural networks is that they only consider the context of the focused word instead of the whole sentence. Recurrent neural networks (RNNs) (Mikolov et al., 2010) can keep the future and the past data persistent by using memory cells with loops in them. However, RNNs are biased towards the most recent data in practice. Long Short Term Memory Networks (LSTMs), special RNNs with Long Short Term memory cells (Graves and Schmidhuber, 2005), are designed to combat the bias problem. LSTMs learn long dependencies in a sequence with the help of Gates (Graves and Schmidhuber, 2005). Gates control how much of the input information to pass to the next LSTM cell, and how much of the previous information to forget.

Given a sentence with n words each of which is represented as a dense vector x_i , a LSTM makes use of $x_{1:n}$ to compute a forward representation \vec{h}_i for the i th word. In general, computing a backward representation \overleftarrow{h}_i would be useful for sequence tagging. Bidirectional LSTM (BiLSTM) is an extension to LSTM which takes into account both the past elements and the future elements in a sequence. A BiLSTM generates both \vec{h}_i and \overleftarrow{h}_i for the i th word in a sequence. The hidden embedding h_i is the concatenation of \vec{h}_i and \overleftarrow{h}_i . BiLSTM mapping is defined as $BiLSTM_\theta: h_i = BiLSTM_\theta(x_{1:n}, i)$ where θ represents trainable parameters in BiLSTM.

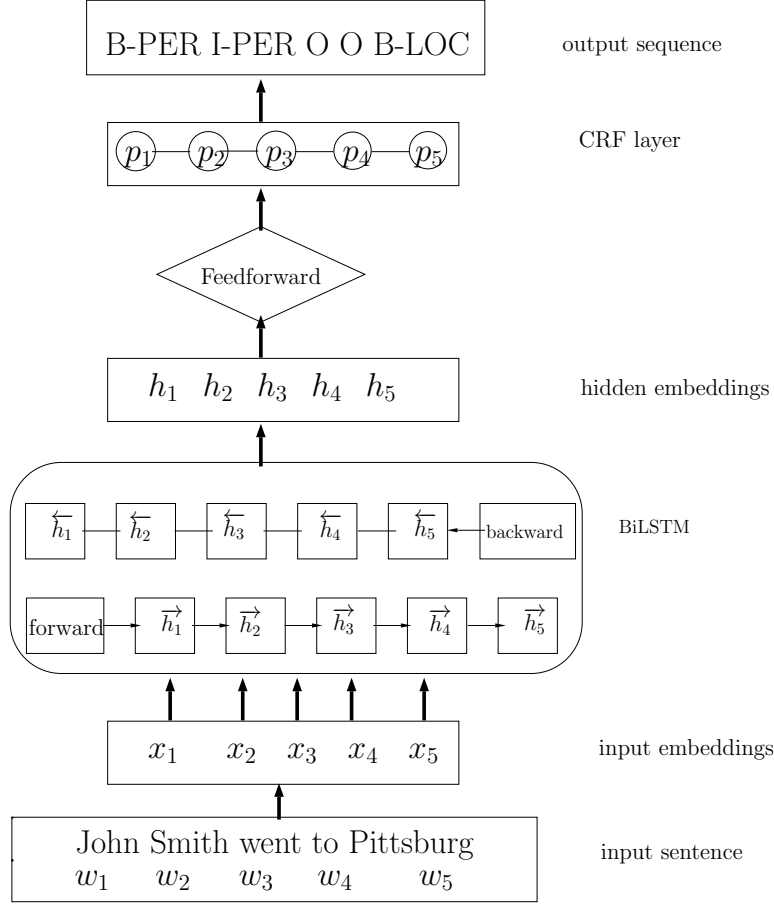


Figure 3.1: The architecture of an NER tagging system using BiLSTM with CRF

We can combine BiLSTM with CRF to make use of sentence level tag information. As in the Feedforward-CRF model, we introduce a transition matrix T to keep the transition score from y_i to y_{i+1} . Given an output sequence Y , the score of the output sequence is given by:

$$S(X|Y) = \sum_i^n T_{i,i+1} + \sum_i^n P_i \quad (3.1)$$

We can use the dynamic programming algorithm to compute the transition matrix and find the optimal output tag sequence. Figure 3.1 illustrates the BiLSTM model with CRF on an NER example.

3.1.2 Character Embeddings

Instead of using hand-engineered features listed in Chapter 2 (like prefixes and the suffixes of a word), we can use a BiLSTM network to construct word representations from characters in it (Lample et al., 2016). This architecture is denoted as Character Embeddings. It has been shown that learning character embeddings has been found useful for capturing morphological evidence (Ling et al., 2015). Figure 3.2 describes the architecture using character embeddings and BiLSTM to generate word embedding for word "Smith". The input to the BiLSTM is the letter sequence of a word. We define a character dictionary which maps each character to a d -dimensional vector representation. The English character dictionary contains uppercase and lowercase letters, numbers, and punctuation. We look up each c_i of the input letter sequence from the dictionary and get the character embedding vectors $X : \{x_1, x_2, \dots, x_n\}$. Then, character embedding vectors X is fed into BiLSTM to generate forward and backward hidden embeddings of the character sequence. We concatenate the last forward hidden embedding \vec{h}_n and the last backward hidden embedding \overleftarrow{h}_1 with the embeddings from word embedding dictionary lookup to obtain the final word embedding.

We add Character Embeddings in the sequence tagging system by concatenating the output of Character Embeddings and the word embeddings from lookup to form input embeddings. Then, we feed the input embeddings to BiLSTM and CRF, and we get the fully structured BiLSTM model (BiLSTM-Char-CRF) for sequence tagging. The character embeddings and word embeddings are learned together during training. There are existing implementations of BiLSTM-Char-CRF, such as NeuroNet by Dernoncourt et al. (2017) which achieves state-of-the-art performance. In order to compare BiLSTM-Char-CRF with other models in this thesis, we re-implement the model.

3.2 Experiments and Results

To evaluate BiLSTM models, we run BiLSTM with three configurations on POS and NER: BiLSTM with word features only (BiLSTM); BiLSTM with Character Embeddings (BiLSTM-Char); BiLSTM with Character Embeddings and a CRF layer (BiLSTM-Char-CRF). We report the performance and decoding speed of the models on Penn Treebank data set for POS, and on CoNLL 2003 data set for NER.

As in the implementation for feedforward models, we implement BiLSTM models using Python and the Tensorflow 1.0 library. The hidden layer size in the BiLSTM of Character Embeddings is set to 50, and the hidden layer size in word sequence BiLSTM is set to 100. The rest of the hyperparameters are the same as the ones in the feedforward network experiments. Since dropout training (Hinton et al., 2012) can improve the performance by encouraging the model to depend on both character embeddings and word embeddings, we apply a dropout mask on the input embeddings before the BiLSTM layer. The dropout

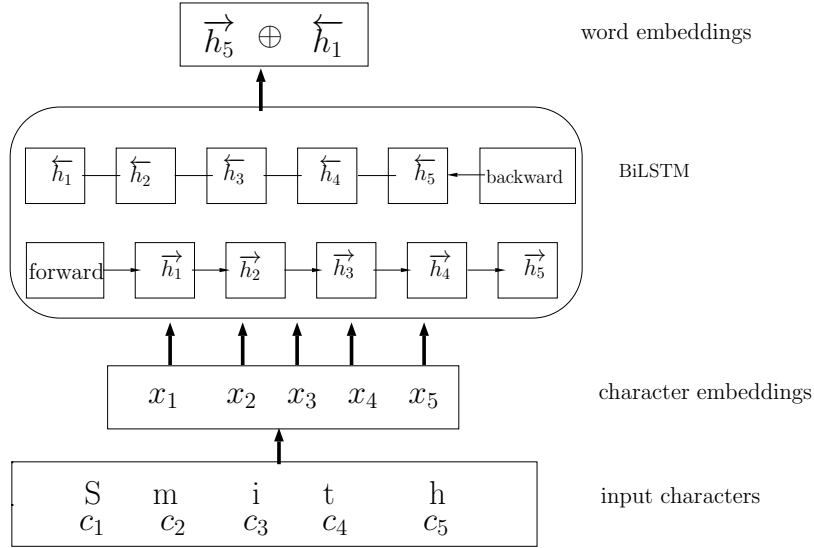


Figure 3.2: The word embedding derived from the character embeddings

rate is set to 0.5 in the experiments. Table 3.1 shows the hyperparameters used in BiLSTM models.

Figure 3.3 illustrates the performance of BiLSTM, BiLSTM-Char, and BiLSTM-Char-CRF. As expected, the fully structured BiLSTM model outperforms the other two. The structure of Character Embeddings increases POS accuracy by 1.17, and increases NER F1 score by 3.34. It reveals that the BiLSTM networks do not heavily depend on features other than words, and Character Embeddings can help improve the performance on NER more than on POS. Compared to BiLSTM-Char, BiLSTM-Char-CRF improves the POS

Table 3.1: Hyperparameters used in BiLSTM Models

Hyperparameters	Values
character embedding size	50
word embedding size	100
feedforward layer size	200
feedforward layer	1
BiLSTM layer size	100
BiLSTM layer	2
optimizer	Adam
learning rate	0.01
batch size	32

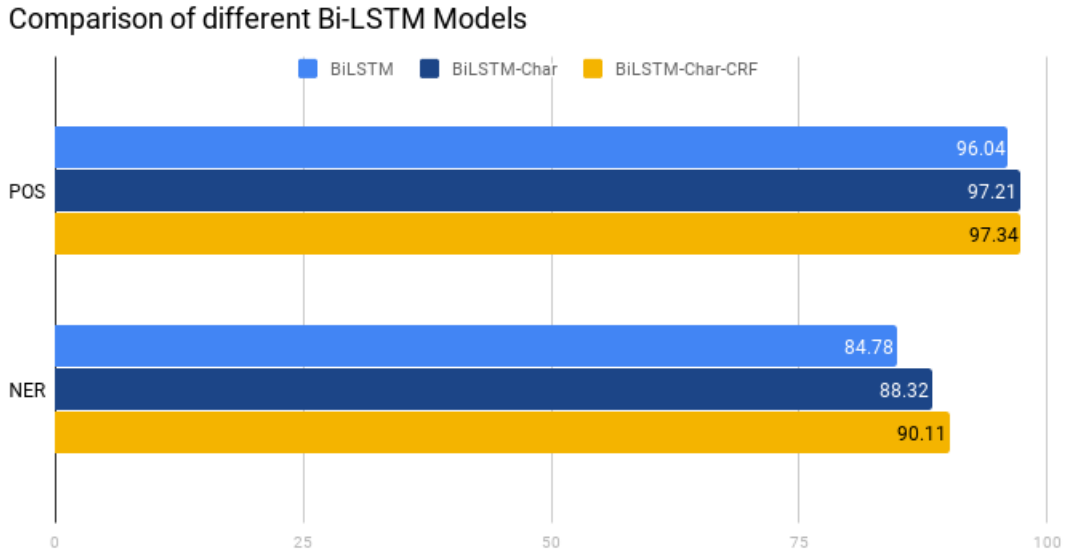


Figure 3.3: Performance comparison between BiLSTM models on POS and NER (accuracy for POS; F1 score for NER)

Table 3.2: BiLSTM Models Accuracy and F1 Score

Model	POS (Accuracy)	NER (F-Score)
BiLSTM	96.01	84.78
BiLSTM-Char	97.21	88.32
BiLSTM-Char-CRF	97.34	90.11

performance by 0.13, and improves the NER performance by 1.79. CRF layer is more helpful for NER for the reason that there are grammar constraints on NER output tag sequence.

Table 3.2 and Table 3.3 describe the final results of the performance and decoding speed of BiLSTM models on the POS and NER. It is obvious that the more features the model has, the slower it would be. Adding a CRF layer will increase the performance but decrease the decoding speed.

Table 3.3: BiLSTM Models Decoding Speed

Model	POS (sentences, words/sec)	NER (sentences, words/sec)
BiLSTM	981(23036)	1637(20740)
BiLSTM-Char	596(13992)	889(11271)
BiLSTM-Char-CRF	383(9009)	795(10100)

Chapter 4

Mention2Vec Models

In this chapter, we introduce two new multi-task models for POS and NER. The multi-task models are based on Mention2Vec which is proposed for NER. We conduct experiments using the multi-task models on POS and NER and compare their performance and decoding speed with Feedforward and BiLSTM-Char-CRF

4.1 Model Description

4.1.1 Feedforward-Mention2Vec for NER

Mention2Vec is a neural network model for NER, which uses BiLSTMs to predict boundaries and entity types separately (Stratos, 2016). We summarize Mention2Vec into three steps. The first step uses a BiLSTM to generate hidden embeddings the same way in the BiLSTM-Char-CRF model. Figure 4.1 illustrates the first step of Mention2Vec on an NER example. We denote the input word sequence as $W : \{w_1, w_2, \dots, w_n\}$, input embeddings as $X : \{x_1, x_2, \dots, x_n\}$ which are the concatenations of the character embeddings and the word embeddings, and the hidden embeddings as $H : \{h_1, h_2, \dots, h_n\}$

The second step of Mention2Vec is boundary detection which predicts $\{I, O, B\}$ labels for words in a sequence. Unlike usual NER labels (B-LOC, I-LOC, ...), boundary labels in this step do not have NER types attached. Mention2Vec takes the hidden embeddings produced in the first step, and feeds them in to a feedforward network to get boundary label probability distribution for each word. Since the boundary labels have strong correlations, Mention2Vec uses CRF to capture correlations and produce boundary label sequences. Figure 4.2 illustrates the second step of Mention2Vec. The output boundary label probability distribution is denoted as p_i for word w_i , and the gold boundary label sequence is denoted as Y_{label} . In each training step, the boundary detection loss is given by the negative log likelihood of the gold boundary label sequence, shown in Equation 4.1

$$L_1(\theta_1) = -\log(p(Y_{label}|h_1, h_2 \dots h_n)) \quad (4.1)$$

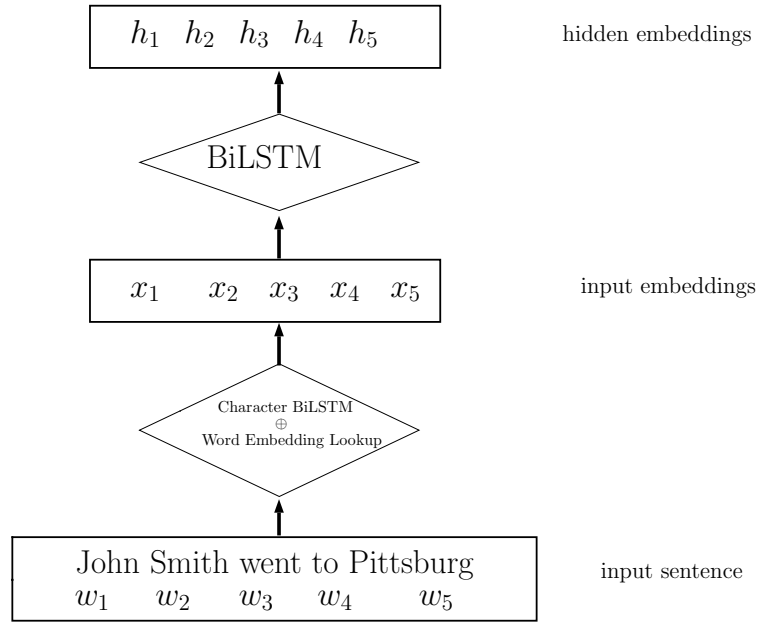


Figure 4.1: The first step of Mention2Vec for NER

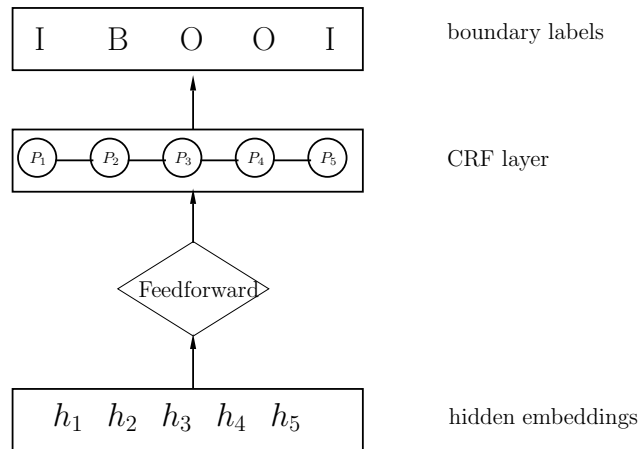


Figure 4.2: The second step of Mention2Vec for NER

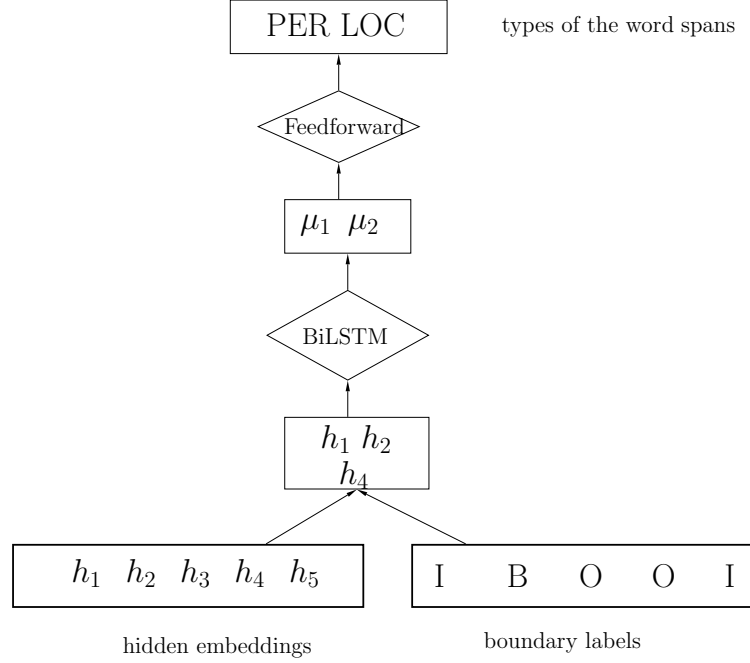


Figure 4.3: The third step of Mention2Vec for NER

The third step of Mention2Vec is type prediction which finds actual types for named entity spans in a sequence. It makes use of the hidden embeddings and the entity boundaries from the previous two steps. The model first looks up the hidden embeddings for the words in entity spans. Then, it feeds the hidden embeddings into an additional BiLSTM and obtains the corresponding vector representations for the word span μ by concatenating the final forward and the backward states produced by the BiLSTM. A feedforward network, at the end, is used to map the word span vector representations μ to type probability distributions. Figure 4.3 illustrates the third step of Mention2Vec.

The gold type output of an input sequence is denoted as Y_{type} . Assuming there are l named entities in a sequence, and the index of the first word in a named entity is represented as s , and the index of the last word in a named entity is represented as e , the type prediction loss is computed by 4.2:

$$L_2(\theta_2) = - \sum_l \log P(r^l | h_s^l \dots h_e^l) \quad (4.2)$$

During training, the model uses the gold boundaries and gold types to compute the boundary detection loss and the type prediction loss. In each training step, the boundary detection loss and the type prediction loss are minimized jointly: the training objective is to find boundary sequences and type sequences that minimize the sum of L_1 and L_2 .

In order to further speed up Mention2Vec as well as capture the correlation between boundary tags, we consider using a different network for boundary detection. Shown in

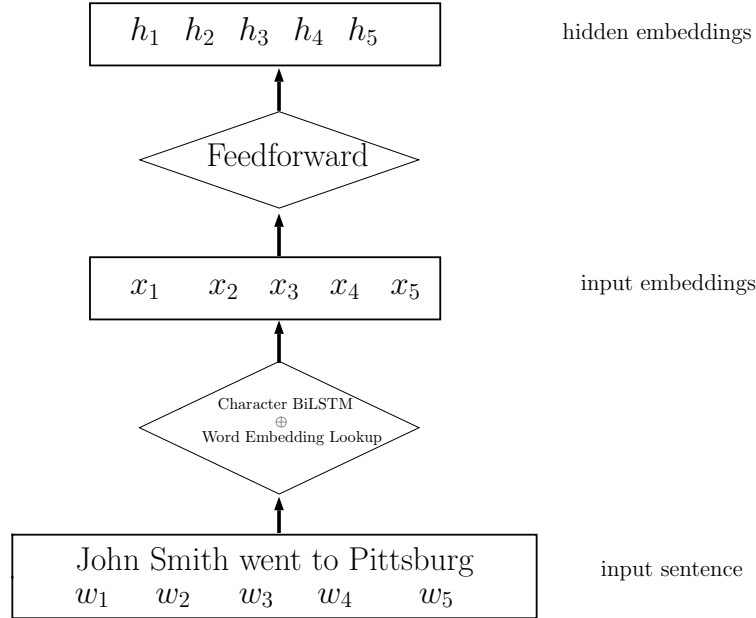


Figure 4.4: The first step of Feedforward-Mention2Vec for NER

Chapter 2 and Chapter 3, Feedforward-CRF produces relatively good performance on NER with faster speed than BiLSTM. We then replace the BiLSTM in the first step of Mention2Vec with a feedforward network. We denote this new model for NER as Feedforward-Mention2Vec. Feedforward-Mention2Vec still has three steps where the second step and the third step are the same as the ones in Mention2Vec. The first step uses a feedforward network to produce the hidden embeddings, which is illustrated in Figure 4.4.

The run time of decoding boundaries is constant in both Feedforward-Mention2Vec and Mention2Vec, since there are only three types of boundaries ('I', 'O', 'B') to decode in the CRF layer. In other CRF based models, like BiLSTM-Char-CRF, the decoding time grows quadratically in the number of named entity types.

4.1.2 BPE-Mention2Vec for POS

In POS, each word in the input sentence is assigned a unique part-of-speech tag. Since there is no part-of-speech span existing in POS, it's not necessary to use a multi-task model on POS. However, we propose a way to convert POS into an NER-like task through the help of Byte Pair Encoding. Inspired by the successful results obtained by using BPE in machine translation (Sennrich et al., 2015), we initially wanted to use BPE to capture morphological decomposition of the words and replace spelling features like prefixes and suffixes. BPE is a compression algorithm which replaces frequent pairs with an unused byte. Sennrich et al.

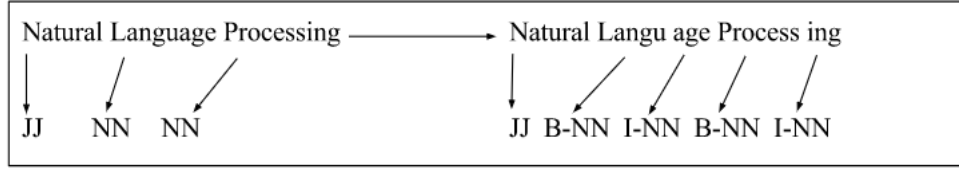


Figure 4.5: An example of using BPE for word segmentation

(2015) presents a way to adapt BPE for word segmentation: using BPE to segment words into subword units of different length, and building a vocabulary dictionary using word frequency. In POS tagging system, we first learn BPE merge operations on the training data. To segment training data and test data, we first split each word in characters and then apply BPE to merge characters into larger chunks. In order to restore the words, we use the "IOB" label scheme to label the subword units. Since there can be multiple subword units sharing the same tag, POS becomes a task similar with NER. Figure 4.5 shows an example of using BPE to segment a sequence with POS tags.

In our proposed BPE-Mention2Vec model for POS, there are three main steps. In the first step, given a sequence, we use BPE to segment the words and convert the corresponding tags using the "IOB" label scheme. After the first step, we have an NER-like task with known boundary of each entity span, so we can apply the same method in Feedforward-Mention2Vec to predict the POS types for entity spans. In the second step, we use a feedforward network to produce the hidden embeddings of the input words. CRF is not used here because the boundary labels are known in both training and decoding. The third step of BPE-Mention2Vec uses a BiLSTM to predict the POS tags for each entity span based on the hidden embeddings and the known boundaries. Figure 4.6 describes the process of using BPE-Mention2Vec to find POS tags for a sequence.

4.2 Experiments and Results

We empirically evaluate the Mention2Vec model and the Feedforward-Mention2Vec model for NER, and the BPE-Mention2Vec model for POS. We implement these models in Python and Tensorflow 1.0. The experiments in this section share the same set of hyperparameters as the feedforward and BiLSTM experiments.

In NER experiments, we compare Feedforward-Mention2Vec with Mention2Vec on the CoNLL 2003 data set. We use the BiLSTM-Char-CRF and Feedforward as baseline models because BiLSTM-Char-CRF achieves the highest F1 score and Feedforward has the fastest decoding speed among the previous models we built. Table 4.1 shows the NER performance and decoding speed of these neural network models on the CoNLL 2003 data set. Mention2Vec achieves 89.51 F1 score which is lower than the 90.11 F1 score of BiLSTM-CRF,

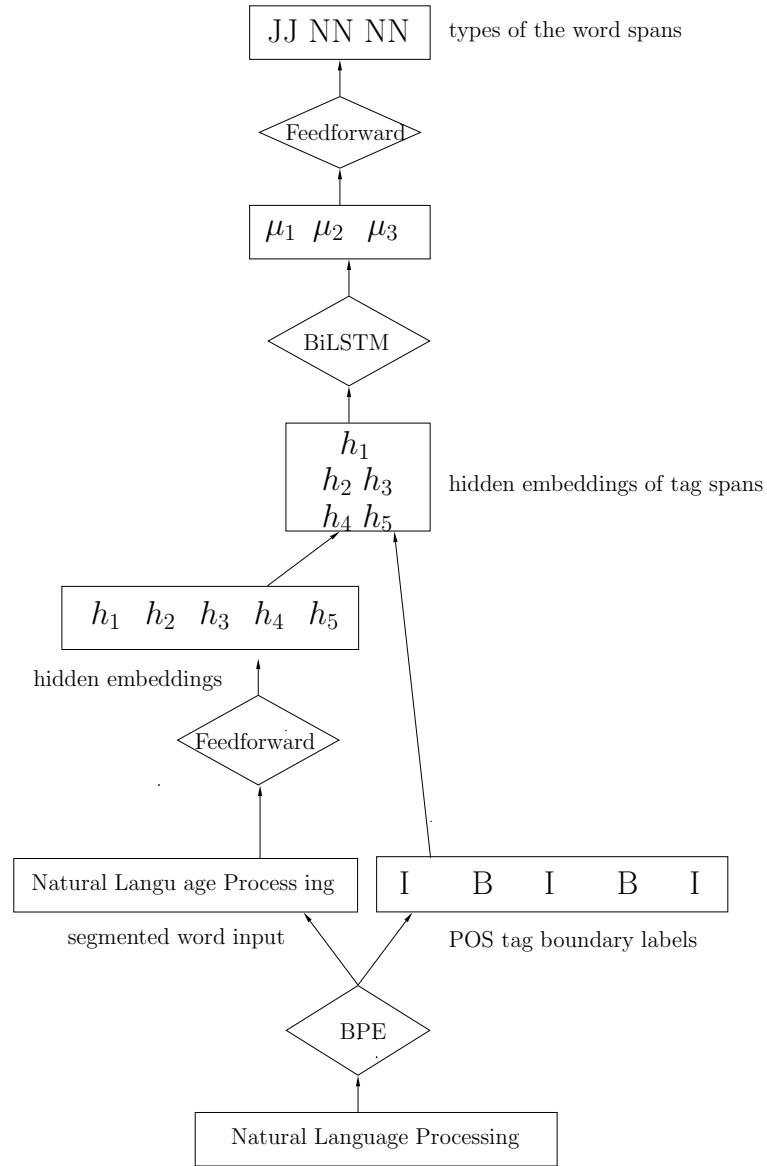


Figure 4.6: An example of using BPE-Mention2Vec to find POS tags

Table 4.1: F1 Score and Decoding Speed Comparison on CoNLL 2003

Model	F1	Speed(sentences, words/sec)
Feedforward	84.12	2117(26819)
BiLSTM-Char-CRF	90.11	795(10100)
Mention2Vec	89.51	765(9701)
Feedforward-Mention2Vec	88.97	1059(13445)

Table 4.2: F1 Score and Decoding Speed Comparison on OntoNotes

Model	F1	Speed(sentences, words/sec)
Feedforward	78.06	1435(22829)
BiLSTM-Char-CRF	86.24	481(7667)
Mention2Vec	85.31	530(8433)
Feedforward-Mention2Vec	83.18	658(10812)

and it is slightly slower than BiLSTM-Char-CRF. Feedforward-Mention2Vec achieves 88.79 F1 score compared to 89.51 of Mention2Vec and 90.11 of BiLSTM-Char-CRF, but it is 1.4 times faster than Mention2Vec and 1.3 times faster than BiLSTM-Char-CRF. The empirical results demonstrate that the Feedforward-Mention2Vec model performs competitively on the CoNLL 2003 NER tagging task while being faster than original Mention2Vec model and the state-of-the-art model. Feedforward is still the fastest but it has the lowest F1 score.

We have noted that Feedforward-Mention2Vec and Mention2Vec scale linearly in the number of named entity types while BiLSTM-Char-CRF grows quadratically. We want to show the time difference of the same model on NER data sets with different numbers of named entity types. We have obtained the performance and decoding speed of these models on CoNLL 2003 which has 4 types of named entity. Then we conduct the experiments on the OntoNotes English data set which has 18 types of named entity. Table 4.2 reports the final results of the experiments. The same model obtains lower F1 score on OntoNotes. Table 4.3 shows the per label performance of Feedforward-Mention2Vec, BiLSTM-Char-CRF, and Feedforward. They reveal that some labels in OntoNotes are classified poorly, such as "WORK OF ART" and "LANGUAGE". We suspect this is due to OntoNotes is more noisy than CoNLL 2003 as OntoNotes data is extracted from a wide variety of sources with more named entity types. In terms of decoding speed, while Feedforward-Mention2Vec is 1.3 times faster than the BiLSTM-Char-CRF model on the data set with 4 named entity types, Feedforward-Mention2Vec is 1.4 times faster on the data set with 18 named entity types. Mention2Vec is slightly slower than BiLSTM-Char-CRF on CoNLL 2003, but it's faster than BiLSTM-Char-CRF on OntoNotes.

In the POS experiments, we compare BPE-Mention2Vec with BiLSTM-Char-CRF which is the state-of-the-art model for POS, and with Feedforward which is the fastest among the models we have built. Table 4.4 demonstrates the performance and decoding speed

Table 4.3: Per-label F1 Score on OntoNotes

Labels	Feedforward	Feedforward-Mention2Vec	BiLSTM-Char-CRF
CARDINAL	72.38	78.03	80.33
DATE:	72.45	80.61	81.24
EVENT:	30.01	45.81	61.40
FAC:	29.08	46.27	57.94
GPE:	89.62	92.84	94.84
LANGUAGE:	45.16	43.42	51.43
LAW:	20.02	43.62	57.97
LOC:	52.79	64.71	73.51
MONEY:	77.73	81.66	86.98
NORP:	84.11	87.80	92.15
ORDINAL:	71.26	72.77	77.24
ORG:	72.05	79.92	85.27
PERCENT:	84.43	89.27	88.73
PERSON:	82.72	86.98	90.65
PRODUCT:	51.66	50.00	64.38
QUANTITY:	65.78	74.77	81.13
TIME:	41.68	56.44	58.85
WORK OF ART:	26.29	37.63	47.21

Table 4.4: POS Tagging System Accuracy and Decoding Speed

Model	Accuracy	Speed(sentences, words/sec)
BILSTM-Char-CRF	97.34	383(9009)
BPE-Mention2Vec	96.04	209(4923)
Feedforward	95.89	1319(30967)

of them on the Penn Treebank data set. BPE-Mention2Vec obtains lower accuracy than BiLSTM-Char-CRF and higher accuracy than Feedforward. Since using word segmentation increases the number of words to be processed and introduces more preprocessing time, BPE-Mention2Vec is slower than the BiLSTM-Char-CRF and Feedforward. The empirical results conclude that BiLSTM-Char-CRF is a better model than BPE-Mention2Vec on POS.

Chapter 5

Discussion

In this Chapter, we put together the final experiment results of the network models presented in this thesis, and provide an analysis of the results.

Table 5.1 records the performance of the neural network models for POS and NER. The POS performance is measured on the Penn Treebank data set, and the NER performance is measured on the CoNLL 2003 data set. It also includes the results from the state-of-the-art models on Penn Treebank and CoNLL 2003. While our implementations obtain slightly lower accuracy score and F1 score than the state-of-art results, we emphasize that our main goal of this thesis is to compare different neural networks, and to present two new multi-task models.

Table 5.2 records the average decoding speed of different neural network models on the Penn Treebank and CoNLL 2003 test data. All experiments are conducted on a GPU. It is obvious that the fewer features used in the same model the faster the decoding speed of the model will be. In general, the greedy tagging systems using feedforward models are faster than the systems using BiLSTM models. Since the CRF based models introduce a transition matrix and uses dynamic programming algorithm to decode the sequence, they are slower than the models without the CRF layer.

Figure 5.1 illustrates the trade-off between performance and decoding speed in POS systems using different neural network models. Among the neural network models for POS, BiLSTM-Char-CRF achieves the best per word accuracy 97.34%, and Feedforward-CRF obtains the second best per word accuracy 97.30%. Feedforward achieves the fastest decoding speed since it employs a simple neural network without extra features: the POS greedy system using the Feedforward model decodes 1319 sentences per second. The line in Figure 5.1 connects BiLSTM-Char-CRF which is the most accurate model and Feedforward which is the fastest model. The models above the line are faster but perform slightly worse than BiLSTM-Char-CRF, such as Feedforward-CRF and BiLSTM-Char. Models below the line such as BPE-Mention2Vec are slower and less accurate, which makes them less ideal

Table 5.1: Neural Network Models Accuracy and F1 Score

Model	Penn Treebank (Accuracy)	CoNLL 2003 (F1 Score)
Ling et al. (2015)	97.78	—
Lample et al. (2016) with CRF	—	90.94
Feedforward	95.89	84.12
Feedforward-History	97.28	86.54
Feedforward-CRF	97.30	87.85
BiLSTM	96.04	84.78
BiLSTM-Char	97.21	88.32
BiLSTM-Char-CRF	97.34	90.11
Feedforward-Mention2Vec	—	88.97
BPE-Mention2Vec	96.04	—

Table 5.2: Neural Network Models Decoding Speed

Model	Penn Treebank (words/sec)	CoNLL 2003 (words/sec)
Feedforward	30967	26819
Feedforward-History	19474	17609
Feedforward-CRF	17877	17412
BiLSTM	23036	20740
BiLSTM-Char	13992	11271
BiLSTM-Char-CRF	9009	10100
Feedforward-Mention2Vec	—	13445
BPE-Mention2Vec	4923	—

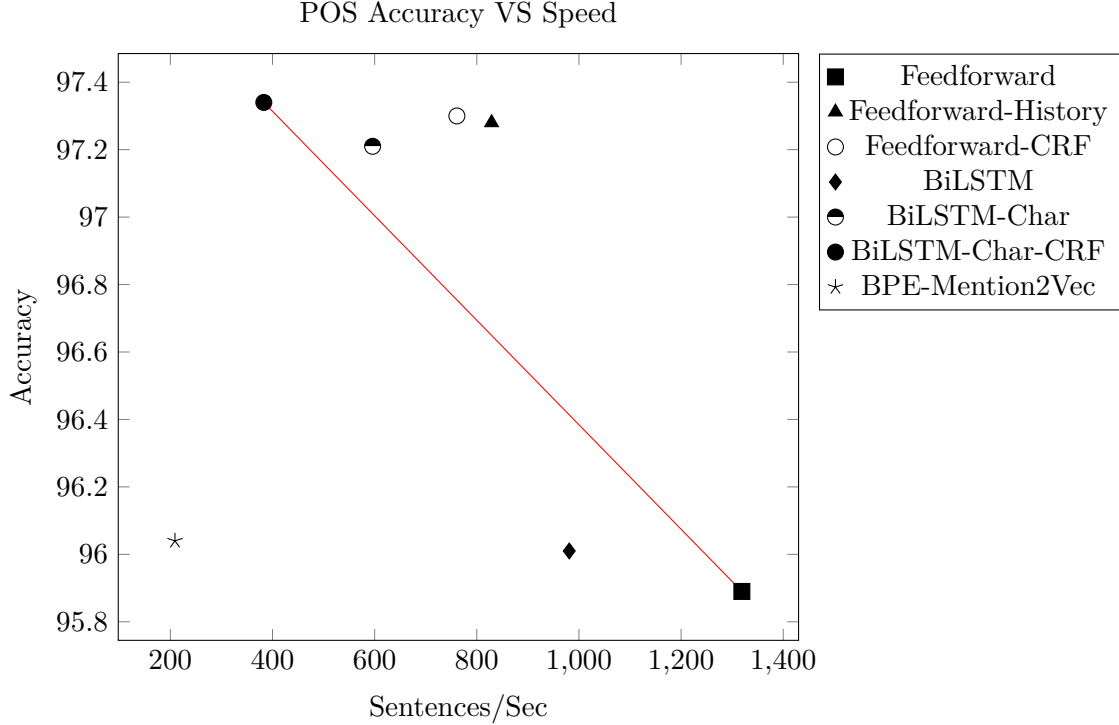


Figure 5.1: Results of the POS system using different Neural Network Models

for POS. Feedforward-History is the fastest model with competitive performance on POS, and it is about 2 times faster than BiLSTM-Char-CRF.

Figure 5.2 illustrates the trade-off between performance and speed of using different neural network models on CoNLL 2003. BiLSTM-Char-CRF achieves the highest F1 Score 90.11, and Feedforward-Mention2Vec obtains the second best F1 Score 88.79. Feedforward achieves the fastest decoding speed since it employs a simple neural network with no extra features: the NER greedy system using the Feedforward model decodes 2177 sentences per second. The line in Figure 5.2 connects the BiLSTM-Char-CRF model which is the most accurate model and the Feedforward model which is the fastest model. The models above the line are faster but perform slightly worse than BiLSTM-Char-CRF, such as Feedforward-CRF and Feedforward-Mention2Vec. Models below the line are slower and less accurate, which makes them less ideal for NER. Feedforward-Mention2Vec model is the fastest model with competitive performance. Feedforward-Mention2Vec obtains 88.79 F1 score which is close to the best performance, and it is 1.3 times faster than the fully structured BiLSTM model.

As illustrated in Figure 5.1 and Figure 5.2, the greedy sequence tagging systems using feedforward networks can achieve comparable performance and faster speed than the systems using recurrent models on POS; the multi-task model (Feedforward-Mention2Vec) performs competitively with the state-of-the-art model (BiLSTM-Char-CRF) on NER.

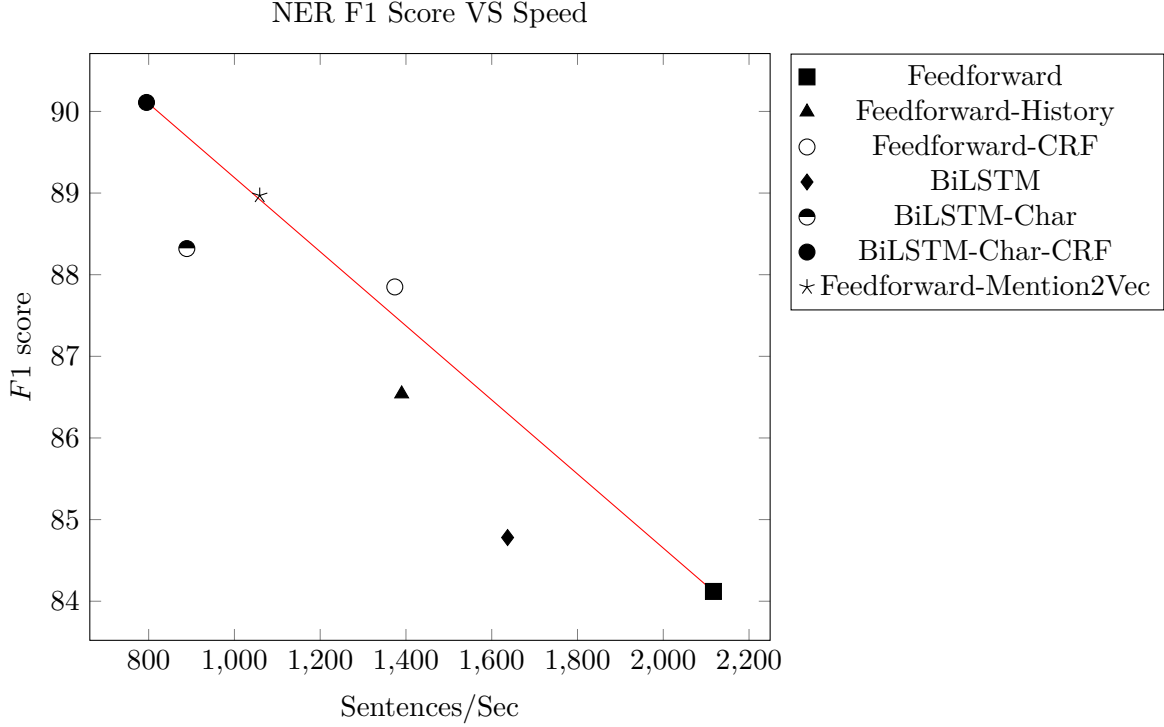


Figure 5.2: Results of the NER system using different Neural Network Models on CoNLL 2003

Table 5.3: F1 Scores and Decoding Speed on CoNLL 2003 and OntoNotes

	F1 Score		Decoding Speed (Words/Sec)	
	CoNLL 2003	OntoNotes	CoNLL 2003	OntoNotes
Feedforward-Mention2Vec	88.97	86.24	13445	10812
BiLSTM-Char-CRF	90.11	83.18	10100	7667

Table 5.3 compares the performance and decoding speed of BiLSTM-Char-CRF and Feedforward-Mention2Vec on CoNLL 2003 and OntoNotes. It shows that Feedforward-Mention2Vec with a simple architecture can achieve competitive performance with BiLSTM-Char-CRF on CoNLL 2003 and OntoNotes. The speed gap between Feedforward-Mention2Vec and BiLSTM-Char-CRF grows as the number of named entity type increases. There is more speed benefit in using a multi-task model like Feedforward-Mention2Vec when there are more named entity types to classify.

Chapter 6

Conclusion and Future Work

This thesis presents and compares different neural network models for sequence tagging tasks. The empirical results reveal that simple feedforward networks can achieve competitive results while being significantly faster than the BiLSTM networks on POS. The empirical results also demonstrate that Feedforward-Mention2Vec performs well on NER, and it is more scalable in the number of named entity types.

6.1 Contribution

In this thesis, we first built three feedforward neural network models for POS and NER: Feedforward, Feedforward-History, and Feedforward-CRF. We conducted experiments to compare the decoding speed and performance between them. Experimental evaluations show that the feedforward models are not strongly dependent on hand engineered features, and the models are able to automatically learn the useful features for making decisions. Compared to Feedforward-CRF, Feedforward-History achieves a similar per-word accuracy score on POS, but it is less accurate on NER.

Then, we re-implemented the state-of-the-art model for sequence tagging: BiLSTM-Char-CRF. We ran the BiLSTM model with different configurations on POS and NER, and we compared the decoding speed and performance between BiLSTM models and feedforward models. As the experiments show, BiLSTM models are robust and are more accurate than feedforward models in general. Since feedforward models have a simpler architecture and fewer parameters, feedforward models are faster than BiLSTM models.

In addition to feedforward models and BiLSTM models, we presented Feedforward-Mention2Vec for NER which is a combination of Feedforward-CRF and Mention2Vec, and BPE-Mention2Vec for POS which is based on BPE and Mention2Vec. Both of these two models are multitasking models. Feedforward-Mention2Vec predicted boundaries of named entity labels and types of named entity labels separately using a feedforward network, a BiLSTM network, and CRF. BPE-Mention2Vec first segmented words using BPE, and pre-

dicted the part-of-speech tags for the subword units using a BiLSTM network. Feedforward-Mention2Vec performs slightly worse than the state-of-the-art model (BiLSTM-Char-CRF) on CoNLL 2003, but its decoding speed is faster than BiLSTM-Char-CRF. As the number of named entity types grows, the decoding speed of Feedforward-Mention2Vec grows linearly while the decoding speed of BiLSTM-Char-CRF grows quadratically. For example, Feedforward-Mention2Vec is 1.3 times faster than BiLSTM-Char-CRF on CoNLL 2003 with 4 named entity types while it is 1.4 times faster on OntoNotes with 18 named entity types.

Lastly, we summarized all the experiment results and compared the performance and decoding speed of all the models presented in this thesis. We provided analysis on the speed and accuracy trade-off in different models. Feedforward is the fastest among all models, since it includes no extra features and employs a simple network architecture. Our re-implementation of BiLSTM-Char-CRF achieves near state-of-the-art performance on POS and NER. Feedforward-History performs better on POS than on NER with a faster decoding speed than BiLSTM-Char-CRF. Feedforward-Mention2Vec performs slightly worse than BiLSTM-Char-CRF on NER, but its decoding speed is faster and grows linearly in the number of named entity types.

6.2 Future Work

There are several potential extensions of this thesis we would like to work on in the future:

First, Shallow Parsing, often referred as Chunking, is another classic sequence tagging task which can be solved using neural network models proposed in this thesis. We would like to compare different neural network models on Shallow Parsing and find the speed and accuracy relationship between these models.

Second, we would like to explore multitasking models which combine POS and NER. We can build a neural network tagging system to train and predict POS tags and NER tags jointly. The intermediate representations would benefit from considering linguistic hierarchies in the training process.

Third, in order to improve the performance of our sequence tagging system, we would like to combine the neural network models presented in this thesis. The easiest way to combine them is to use a voting method in post-process. During tagging, each model will predict a label for the input word, and we will take a majority vote on the output of the three models. The input word will be tagged with the label with the highest vote.

Bibliography

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. (2016). Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283.
- Alberti, C., Andor, D., Bogatyy, I., Collins, M., Gillick, D., Kong, L., Koo, T., Ma, J., Omernick, M., Petrov, S., et al. (2017). Syntaxnet models for the conll 2017 shared task. *arXiv preprint arXiv:1703.04929*.
- Andor, D., Alberti, C., Weiss, D., Severyn, A., Presta, A., Ganchev, K., Petrov, S., and Collins, M. (2016). Globally normalized transition-based neural networks. *arXiv preprint arXiv:1603.06042*.
- Chen, D. and Manning, C. D. (2014). A fast and accurate dependency parser using neural networks. In *Emnlp*, pages 740–750.
- Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014). On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*.
- Collins, M. (2002). Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*, pages 1–8. Association for Computational Linguistics.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537.
- Dernoncourt, F., Lee, J. Y., and Szolovits, P. (2017). NeuroNER: an easy-to-use program for named-entity recognition based on neural networks. *Conference on Empirical Methods on Natural Language Processing (EMNLP)*.
- Dugas, C., Bengio, Y., Bélisle, F., Nadeau, C., and Garcia, R. (2001). Incorporating second-order functional knowledge for better option pricing. In *Advances in neural information processing systems*, pages 472–478.

- Dyer, C., Ballesteros, M., Ling, W., Matthews, A., and Smith, N. A. (2015). Transition-based dependency parsing with stack long short-term memory. *arXiv preprint arXiv:1505.08075*.
- Florian, R., Ittycheriah, A., Jing, H., and Zhang, T. (2003). Named entity recognition through classifier combination. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 168–171. Association for Computational Linguistics.
- Gage, P. (1994). A new algorithm for data compression. *The C Users Journal*, 12(2):23–38.
- Glorot, X., Bordes, A., and Bengio, Y. (2011). Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 513–520.
- Graves, A., Mohamed, A.-r., and Hinton, G. (2013). Speech recognition with deep recurrent neural networks. In *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*, pages 6645–6649. IEEE.
- Graves, A. and Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5):602–610.
- Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory.
- Hovy, E., Marcus, M., Palmer, M., Ramshaw, L., and Weischedel, R. (2006). Ontonotes: the 90% solution. In *Proceedings of the human language technology conference of the NAACL, Companion Volume: Short Papers*, pages 57–60. Association for Computational Linguistics.
- Huang, Z., Xu, W., and Yu, K. (2015). Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.
- Kingma, D. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Lafferty, J., McCallum, A., and Pereira, F. C. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data.
- Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., and Dyer, C. (2016). Neural architectures for named entity recognition. *arXiv preprint arXiv:1603.01360*.

- Ling, W., Luís, T., Marujo, L., Astudillo, R. F., Amir, S., Dyer, C., Black, A. W., and Trancoso, I. (2015). Finding function in form: Compositional character models for open vocabulary word representation. *arXiv preprint arXiv:1508.02096*.
- Marcus, M. P., Marcinkiewicz, M. A., and Santorini, B. (1993). Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330.
- McCallum, A., Freitag, D., and Pereira, F. C. (2000). Maximum entropy markov models for information extraction and segmentation. In *Icml*, volume 17, pages 591–598.
- Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., and Khudanpur, S. (2010). Recurrent neural network based language model. In *Interspeech*, volume 2, page 3.
- Nair, V. and Hinton, G. E. (2010). Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814.
- Nivre, J. and Scholz, M. (2004). Deterministic dependency parsing of english text. In *Proceedings of the 20th international conference on Computational Linguistics*, page 64. Association for Computational Linguistics.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *EMNLP*, volume 14, pages 1532–1543.
- Ratnaparkhi, A. et al. (1996). A maximum entropy model for part-of-speech tagging. In *Proceedings of the conference on empirical methods in natural language processing*, volume 1, pages 133–142. Philadelphia, PA.
- Sennrich, R., Haddow, B., and Birch, A. (2015). Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.
- Stratos, K. (2016). Mention2vec: Entity identification as multitasking. *arXiv preprint arXiv:1612.02706*.
- Tjong Kim Sang, E. F. and De Meulder, F. (2003). Introduction to the conll-2003 shared task: Language-independent named entity recognition. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 142–147. Association for Computational Linguistics.
- Toutanova, K., Klein, D., Manning, C. D., and Singer, Y. (2003). Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*, pages 173–180. Association for Computational Linguistics.