# Neural Networks for Linguistic Structured Prediction and Their Interpretability

Xuezhe Ma

CMU-LTI-20-001

Language Technologies Institute
School of Computer Science
Carnegie Mellon University
5000 Forbes Ave., Pittsburgh, PA 15213
www.lti.cs.cmu.edu

**Thesis Committee:**
Eduard Hovy (Chair), Carnegie Mellon University
Jaime Carbonell, Carnegie Mellon University
Yulia Tsvetkov, Carnegie Mellon University
Graham Neubig, Carnegie Mellon University
Joakim Nivre, Uppsala University

*Submitted in partial fulfillment of the requirements*
*for the degree of Doctor of Philosophy*
*n Language and Information Technologies*

# Abstract

Linguistic structured prediction, such as sequence labeling, syntactic and semantic parsing, and coreference resolution, is one of the first stages in deep language understanding and its importance has been well recognized in the natural language processing community, and has been applied to a wide range of down-stream tasks.

Most traditional high performance linguistic structured prediction models are linear statistical models, including Hidden Markov Models (HMM) and Conditional Random Fields (CRF), which rely heavily on hand-crafted features and task-specific resources. However, such task-specific knowledge is costly to develop, making structured prediction models difficult to adapt to new tasks or new domains. In the past few years, non-linear neural networks with as input distributed word representations have been broadly applied to NLP problems with great success. By utilizing distributed representations as inputs, these systems are capable of learning hidden representations directly from data instead of manually designing hand-crafted features.

Despite the impressive empirical successes of applying neural networks to linguistic structured prediction tasks, there are at least two major problems: 1) there is no a consistent architecture for, at least of components of, different structured prediction tasks that is able to be trained in a truely end-to-end setting. 2) The end-to-end training paradigm, however, comes at the expense of model interpretability: understanding the role of different parts of the deep neural network is difficult.

In this thesis, we will discuss the two of the major problems in current neural models, and attempt to provide solutions to address them. In the first part of this thesis, we introduce a consistent neural architecture for the encoding component, named BLSTM-CNNs, across different structured prediction tasks. It is a truely end-to-end model requiring no task-specific resources, feature engineering, or data pre-processing beyond pre-trained word embeddings on unlabeled corpora. Thus, our model can be easily applied to a wide range of structured prediction tasks on different languages and domains. We apply this encoding architecture to different tasks including sequence labeling and graph and transition-based dependency parsing, combined with different

structured output layers, achieving state-of-the-art performance.

In the second part of this thesis, we use probing methods to investigate learning properties of deep neural networks with dependency parsing as a test bed. We first apply probes to neural dependency parsing models and demonstrate that using probes with different expressiveness leads to inconsistent observations. Based on our findings, we propose to interpret performance of probing tasks with two separate metrics, *capacity* and *accessibility*, which are associated with probe expressiveness. Specifically, *capacity* measures how much information has been encoded, while *accessibility* measures how easily the information can be detected. Then, we conduct systematic experiments to illustrate two learning properties of deep neural networks: (i) *laziness* – storing information in a way that requires minimal efforts; (ii) *targetedness* – filtering out from internal representations information that is unnecessary for the target task.

# Acknowledgments

First and foremost, my greatest thanks go to my advisor, Eduard Hovy. Thank you for welcoming me to Edvisees. I'm especially lucky to have found an advisor who gives me freedom to pursue my ideas, and research I am interested in. At the same time, thank you for pushing me back on track from time to time and making sure this PhD comes to an end. I value your broad perspective and long-term vision, but most of all, your kindness and consideration have made this a pleasant journey.

I would like to thank my committee members, Jaime Carbonell, Yulia Tsvetkov, Graham Neubig and Joakim Nivre, for providing helpful feedback and advice. Thank you for bearing with my tardiness. Your comments have made this thesis more complete, and gave me plenty of ideas for future research. First, I want to cherish the memory of Jaime Carbonell, who was a so smart and nice person and I was so fortunate to have him in my thesis committee. Then, I want to show my special appreciation to Yulia Tsvetkov for her urgent help in the special situation. I thank Graham Neubig, with whom I had a lot collaborations that were very enjoyable. I thank Joakim Nivre, , for serving on my committee and always being open to intelligent advice and comments on earlier drafts of this theses.

I realize how fortunate and privileged I am to work in the Edvisees group, a welcoming environment where I could talk to people on anything research and life, hang out for dinner or beer, and simply be happy to come to the office every day. Thanks to past and present members of Edvisees group for mking this such a great environment: Zhengzhong Liu, Sujay Jauhar, Pradeep Dasigi, Diyi Yang, Nicolas Fauceglia, Yiu-Chang Lin, Jun Araki, Dongyeop Kang, Varun Gangal, Jiarui Xu, Xianyang Chen, Xiang Kong, and Zhisong Zhang, as well as all the MEng and undergraduate students I interacted with over the years.

During my PhD., I have done two wonderful internships at Allen Institute for Artificial Intelligence and Petuum Inc. I thank my mentors Waleed Ammar, najmeh sadoughi and Shanghang Zhang when I worked at these places.

Collaboration is a big lesson that I learned, and also a fun part of graduate school.

I thank my fellow collaborators: Zhiting Hu, Zhengzhong Liu, Di Wang, Zihang Dai, Qizhe Xie, Nicolas Fauceglia, Yiu-Chang Lin, Xiang Kong, Zhisong Zhang, Wasi Uddin Ahmad, Chunting Zhou, Yingkai Gao, Yaoliang Yu, Yuntian Deng, Pengcheng Yin, Jingzhou Liu, Zecong Hu, Junjie Hu, and Shanghang Zhang, Diyi Yang, Nanyun Peng and Kai-Wei Chang. In particular, Nanyun and Diyi — they provide me with great professional suggestions for my Ph.D and academic career.

Much appreciation and gratitude to Stacey Young, our unique administrative assistant, for all her assistance and for making our life much easier.

To my parents, thank you for supporting me in so many ways. The foundations you instilled in me since a young age have directly impacted my success. I dedicate this thesis to you, in love and appreciation.

To Chunting, thank you for standing with me every step of the way. It's been a long, and sometimes hard ride, and there were times I didn't think I'll make it to the finish line. Your constant love and support have made this possible. I couldn't have done this without you.

# Contents

# List of Figures

# List of Tables

xviii

# Chapter 1

# Introduction

Teaching machines to understand human language documents is one of the most elusive and long-standing challenges in Artificial Intelligence. Linguistic structured prediction, such as sequence labeling (Lafferty et al., 2001; Ratinov and Roth, 2009; Passos et al., 2014; Luo et al., 2015), syntactic and semantic parsing Nivre and Scholz (2004); McDonald et al. (2005a); Koo and Collins (2010); Ma and Zhao (2012a,b); Chen and Manning (2014); Ma and Hovy (2015), and coreference resolution (Ng, 2010; Durrett and Klein, 2013; Ma et al., 2016), is one of the first stages in deep language understanding and its importance has been well recognized in the natural language processing (NLP) community.

Many problems in machine learning involve structured prediction, i.e., predicting a group of outputs that depend on each other. Many NLP systems, sentiment analysis (Tai et al., 2015), machine translation (Xie et al., 2011; Bastings et al., 2017), information extraction (Nguyen et al., 2009; Angeli et al., 2015; Peng et al., 2017), word sense disambiguation (Fauceglia et al., 2015), and low-resource languages processing (McDonald et al., 2013; Ma and Xia, 2014), are becoming more sophisticated, in part because of utilizing structured knowledge such as part-of-speech (POS) tags, dependency parsing trees, and entity/event coreference information. Figure 1.1 illustrates examples for three classic linguistic structured prediction tasks:

- **part-of-speech tagging:** to assign to each word in a sentence its part-of-speech (POS) tag to represent the syntax function. For example, in the sentence *I saw a girl with a telescope.* in Figure 1.1, *I* is a pronoun, *saw* is a verb is its past tense, *a* is a determiner, *girl* and

1

**Part-of-Speech Tagging**

**Dependency Parsing**

**Named Entity Recognition**

Figure 1.1: Examples for three linguistic structured prediction tasks: Part-of-speech tagging, named entity recognition and dependency parsing.

*telescope* are common nouns, and *with* is a preposition.

- **named entity recognition:** to recognize the spans of named entities in a sentence. For example, *Arsène Wegner* is a person's name, while *Arsenal* is an organization.

- **dependency parsing:** to analyze the syntactic dependency structure of a given sentence to represent syntactic relations between words. Taking the first sentence *I saw a girl with a telescope.* as an example again, *I* is the subject and *a girl* is the object of the verb *saw*, while *with a telescope* as a while is a prepositional phrase which describes a preposition of manner for the verb.

## 1.1   Feature Engineering in Traditional Structured Predictions

As to the classical techniques applied in structured prediction, there are three main streams: i) rule-based approaches; ii) unsupervised learning approaches; and iii) feature-based supervised learning approaches. Among these approaches, feature-based supervised learning approaches, such as Hidden Markov Models (HMM) and Conditional Random Fields (CRF) (Ratinov and Roth, 2009; Passos et al., 2014; Luo et al., 2015), have stood out for the impressive performance across several linguistic structured prediction tasks.

One of the main challenges of these feature-based approaches is that they rely heavily on

hand-crafted features and task-specific resources. For example, English POS taggers benefit from carefully designed word spelling features; orthographic features and external resources such as gazetteers are widely used in NER; carefully selected combinations of head and context words, and their corresponding morphological features, significantly improve the performance of dependency parsers. However, such task-specific knowledge is costly to develop (Ma and Xia, 2014), making structured prediction models difficult to adapt to new tasks or new domains. More discussions on feature representations are in Section 2.4.

## 1.2 Research Objectives

Representation learning techniques based on deep neural networks (LeCun et al., 2015) has re-emerged as one viable solution to this challenge, and fundamentally transformed the conventional feature engineering paradigm. Representation learning, in principle, enables automatic extraction of dense, continuous feature representations for downstream tasks via the end-to-end learning paradigm (Bengio et al., 2013; LeCun et al., 2015). Unlike traditional linear models that rely on manual feature engineering, deep neural models are able to learn complex and intricate features from data via non-linear activation functions. The expressive power of these representations has led to a number of impressive empirical successes in natural language processing (NLP) and other spheres of artificial intelligence (AI), such as computer vision (Krizhevsky et al., 2012), control (Mnih et al., 2015), robotics (Levine et al., 2016) and several sub-fields in machine learning (Goodfellow et al., 2016). In particular, deep learning has revolutionized natural language processing (NLP), with the primary advancement that words, concepts and contexts which were previously represented as a set of sparse discrete features can be represented as dense, real-valued vectors (Mikolov et al., 2013; Pennington et al., 2014; Graves et al., 2013). Compared to discrete representations which are sparse and can only attain coarse relationships, continuous representations can capture fine-grained similarities between objects. Furthermore, these continuous representations can be updated via the end-to-end representation learning paradigm to optimize the entire neural networks towards the ultimate tasks, such as machine translation (Sutskever et al., 2014; Bahdanau et al., 2015), structured prediction (Graves, 2008; Collobert et al., 2011) and

language generation (Mikolov et al., 2010; Graves, 2013).

Various facets of the problem statement underlying this dissertation are addressed in the scientific articles (either published or accepted for publication) that constitute this dissertation. These articles deal with two distinct research questions.

## Question 1: How to design a consistent neural network architecture that can be applied to different structured prediction tasks?

Despite the impressive empirical successes of applying neural networks to linguistic structured prediction tasks, there is no a consistent architecture for, at least components of, different structured prediction tasks. For different tasks, we have to design specific neural architectures. Therefore, in this thesis, we aim to introduce a consistent neural network architectures that not only support truly end-to-end learning of features from task-oriented (labeled) data, but also be applicable to different tasks. Specifically, the first part of the thesis concentrates on introducing a consistent neural architecture to encode a sentence into its continuous feature representations (one vector for each word), and further applying this neural encoder to linguistic structured prediction tasks.

## Question 2: In what way can deep neural structured prediction models memorize and process linguistic information in their internal representations?

From the perspective of model interpretability, understanding the role of different parts of the deep neural network is difficult. The end-to-end training paradigm significantly simplifies the hand-crafted feature engineering process in traditional feature-based machine learning (ML) systems. This, however, often comes at the expense of model interpretability. Unlike traditional feature-engineered NLP systems whose features, e.g. morphological properties, syntactic categories or semantic relations, are more easily understood by humans, it is more difficult to understand what happens in the internal components of an end-to-end neural network. Such deep neural models are sometimes perceived as "black-box", hindering research efforts and limiting their utility to society (Belinkov, 2018). Therefore, the second part of this thesis focuses on the interpretability of neural networks by investigating the learning properties of neural networks that help us understand the behaviors of the representation learning procedure. In particular, in the second part, we try to design experiments to probe the internal representations of neural networks to understand how deep neural structured prediction models memorize and process linguistic information.

## 1.3  Summary of Contributions

The contributions of this thesis are manifold, summarized as follows:

- We introduced a consistent neural architecture for sentence encoding, which is among the first to research end-to-end neural networks on linguistic structured prediction. In particular, this encoder is capable of capturing both word- and character-level information, and is truly end-to-end, entirely getting rid of feature engineering. Based on this encoder, we proposed three neural networks — one for sequence labeling and two for dependency parsing, demonstrating superior performance on various datasets in different languages. It is an important step towards end-to-end representation learning for texts and inspires following works such as deep contextualized representations (Devlin et al., 2019).

- We made the effort to better understand what linguistic knowledge neural dependency parsing models have actually learned, with probing methods. We concluded that the accuracy of probes with different expressiveness does not consistently reflect the quantity of the encoded information, and further explicitly propose to measure the information encoded in representations instead using two complementary metrics, by taking a prediction view of probing accuracy.

- Finally, we pioneered a new research direction in using probing methods to investigate the learning behaviors of deep neural networks. Specifically, we formally define two learning properties of deep neural networks: *laziness* — information, if already been encoded in some components of a neural network, will not be propagated to other components; and *targetedness* — — information that is unnecessary for the ultimate objective will be filtered out, and conduct systematic experiments to illustrate them.

## 1.4  Thesis Outline

Following the two central themes that we just discussed, this thesis consists of two parts — PART I neural architectures and PART II interpretability.

Part I of this thesis focuses on developing neural networks for linguistic structured prediction

tasks. In Chapter 2, we first visit some background of linguistic structured prediction, end-to-end learning paradigm and give an overview of the history and recent development of neural representation learning on linguistic structured prediction. In the last section of this chapter, we propose a consistent neural architecture, named BLSTM-CNNs, for the encoding component (encoder) across different structured prediction tasks. Next, by stacking different structured decoding layers on top of this encoder, we proposed deep neural models for different linguistic structured prediction tasks.

In Chapter 3, we apply BLSTM-CNNs to sequence labeling tasks, by combining a sequential CRF on top of it, to jointly decode labels for the whole sentence. With no feature engineering and data pre-processing, our model surpassed the performance of all prior models, achieving state-of-the-art performance on two classic sequence labeling tasks — part-of-speech (POS) tagging and named entity recognition (NER). This work was published as Ma and Hovy (2016).

In Chapter 4 and 5, we demonstrate applications of LSTM-CNNs to graph-based and transition-based dependency parsers. For graph-based dependency parsing, we proposed the NeuroMST parser, which constructs a probabilistic structured layer on top of BLSTM-CNNs to define the conditional distribution over all dependency trees. Benefiting from the probabilistic structured output layer, we can use negative log-likelihood as the training objective, where the partition function and marginals can be computed via Kirchhoff's Matrix-Tree Theorem (Tutte, 1984; Smith and Smith, 2007). For transition-based dependency parsing, we proposed a novel neural architecture, *stack-pointer networks* (**STACKPTR**). Unlike traditional transition-based parser with left-to-right shift-reduce actions, the STACKPTR parser performs parsing in an incremental, top-down, depth-first fashion. This architecture makes it possible to capture information from the whole sentence and all the previously derived subtrees, while maintaining a number of parsing steps linear in the sentence length. We evaluate our two parsers on 29 treebanks across 20 languages and different dependency annotation schemas, achieving state-of-the-art performance on most of them. This work is presented in Ma and Hovy (2017) and Ma et al. (2018).

Part II of this thesis focuses on how to interpret the learning behaviors of neural networks.

In Chapter 6, we first briefly discuss the terminological issues regarding analysis and interpretation in machine learning. Then we revisit the probing method (Ettinger et al., 2016; Shi et al.,

6

2016; Belinkov, 2018), and apply it to investigate how part-of-speech information are memorized in neural dependency parsers.

In Chapter 7, we conduct systematic experiments to illustrate two learning properties of deep neural networks: (i) *laziness* – modules of a neural network will not actively learn information that is already learned by other modules; (ii) *targetedness* – information, if unnecessary for the end task, will be filtered out from the internal representations.

In Chapter 8, we finally conclude and discuss future work and open questions in this field.

# Part I

# Neural Models for End-to-end Linguistic Structured Prediction

# Chapter 2

# Background

## 2.1 General Framework of Structured Prediction

Supervised computational methods for the linguistic structured prediction of sentences have been at the forefront of natural language processing research since its inception (Rabiner and Juang, 1986; Eisner, 1996; Tjong Kim Sang and Veenstra, 1999; Charniak, 2000; Lafferty et al., 2001; Collins, 2003). Figure 2.1 graphically displays the framework we will assume for a linguistic structured prediction system.

First, a system should define a *feature extractor*, which takes as input raw sentences and outputs representations of each word that are mathematically and computationally convenient to process for machine learning algorithms. Second, the system should have a *learning algorithm* that takes the training data as input to compute the loss function. Then an optimization algorithm associated with the learning algorithm updates the parameters of the system according to the loss function. This process of producing a structured prediction model from a training set is usually called *training* or *learning*. At last, the model consists of a *decoding algorithm*, a.k.a *inference algorithm*, which specifies how to use the model for prediction. That is , when a new sentence is given to the model, the decoding algorithm uses the parameter specifications in the model to produce a structured output.

Figure 2.1: Outline of generic framework of structured prediction.

### 2.1.1 Formal Definition

In the rest of this thesis, we use the following notations: $\mathbf{x} = \{w_1, \ldots, w_n\}$ represents a generic input sentence, where $w_i$ is the $i$th word. $\mathbf{y}$ represents a generic structured output, e.g. a sequence of tags or a dependency tree. $T(\mathbf{x})$ is used to denote the set of all valid structured outputs $\mathbf{y}$ for sentence $\mathbf{x}$. $D = \{(\mathbf{x}_1, \mathbf{y}_1) \ldots, (\mathbf{x}_N, \mathbf{y}_N)\}$ denotes our training sample, where $(\mathbf{x}_i, \mathbf{y}_i), i = 1, \ldots, N$, are usually i.i.d. samples.

### 2.1.2 Feature Extractor

$\phi(\mathbf{x})$ denotes the feature representation of $\mathbf{x}$, output from the feature extractor ($\phi(w_i)$ is the corresponding representation for $w_i$). For traditional feature-based models, $\phi(\mathbf{x})$ is an abstraction over text where a word is represented by one or many Boolean, numeric, or nominal values ($\phi(w_i)$). For deep learning models, $\phi(\mathbf{x})$ is no longer pre-defined vectors, but continuous vectors output from neural networks. More discussions about feature representations are in Section 2.4.

### 2.1.3 Learning Algorithms

This thesis considers probabilistic models for structured prediction, which defines a family of conditional probability $P_\theta(\mathbf{y}|\mathbf{x})$ over all $\mathbf{y}$ given sentence $\mathbf{x}$, where $\theta \in \Theta$ is the set of parameters of this model. In the context of *maximal likelihood estimation* (MLE), parameters $\theta$ is optimized

to minimize the negative log-likelihood:

$$\min_{\theta \in \Theta} L(D; \theta) = \min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^{N} -\log P_\theta(\mathbf{y}_i | \mathbf{x}_i) = \min_{\theta \in \Theta} \mathrm{E}_{\tilde{P}(\mathbf{y}|\mathbf{x})}[-\log P_\theta(\mathbf{y}|\mathbf{x})] \qquad (2.1)$$

where $\tilde{P}(\mathbf{y}|\mathbf{x})$ is the empirical distribution derived from training data $D$. The learning algorithm is to accomplish the computation of $P_\theta(\mathbf{y}|\mathbf{x})$ for each sentence $\mathbf{x}$ and its structured output $\mathbf{y}$, and its gradients for parameter updates:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_\theta L(D; \theta^{(t)}) \qquad (2.2)$$

where $\eta$ is the learning rate of the gradient decent update in (2.2). For example, in conditional random fields (CRF) (Lafferty et al., 2001), the computation of $L(D; \theta)$ and its gradient can be solved using the Viterbi algorithm (Forney, 1973).

Besides MLE with gradient decent, several researchers developed alternative learning algorithms. Collins (2002) presented the averaged perceptron algorithm for discriminative structured learning. One problem with the perceptron algorithm is that it does not optimize any notion of classification margin, which is widely accepted to reduce generalization error (Boser et al., 1992), leading to a several margin-based learning algorithms. In particular, Crammer and Singer (2001) presented a natural approach to large-margin multi-class classification. Taskar et al. (2004) extended it to structured classification by introducing the Maximum Margin Markov networks ($M^3N$) algorithm. The margin infused relaxed algirithm (MIRA) (Crammer and Singer, 2003; Crammer et al., 2006) employs this optimization within a online framework. McDonald et al. (2005a) applied it to the training of dependency parsers. Daumé et al. (2009) proposed the SEARN algorithm for integrating search and learning to solve structured prediction problems.

### 2.1.4   Decoding Algorithms

The decoding algorithm, on the other hand, is to search for the output $\mathbf{y}^*$ with the highest conditional probability:

$$\mathbf{y}^* = \operatorname*{argmax}_{\mathbf{y} \in T(\mathbf{x})} P_\theta(\mathbf{y}|\mathbf{x}) \qquad (2.3)$$

Since the cardinality of $T(\mathbf{x})$ grows exponentially with the length of the sentence $\mathbf{x}$, it is infeasible to perform exhaustive search directly to sovle Eq. (2.3). Observe that both the learning and

| PRP | VBD | DT | NN | IN | DT | NN |
|-----|-----|-----|------|------|-----|-----------|
| I | saw | a | girl | with | a | telescope |

Part-of-Speech Tagging

| B-PER | I-PER | O | O | B-TTL | O | B-ORG | O | B-DATE |
|-------|--------|-----|-------|---------|-----|---------|-----|--------|
| Arsene | Wenger | was | named | manager | of | Arsenal | in | 1996 |

Named Entity Recognition

Figure 2.2: Sequence labeling formulation of Part-of-speech tagging and named entity recognition. For NER, we use the BIO tagging schema.

decoding algorithms rely on the formulation of the probabilistic model $P_\theta(\mathbf{y}|\mathbf{x})$, and for certain formulations, efficient algorithms exist for solving Eq. (2.1) and (2.3).

Unlike probabilistic structured prediction models, some models are trained by online learning algorithms such as averaged structured perceptron (Freund and Schapire, 1999; Collins, 2002) or margin infused relaxed algorithm (Crammer and Singer, 2003; Crammer et al., 2006; McDonald and Pereira, 2006). For these models, only the algorithm for decoding is required.

## 2.2  Sequence Labeling

The task of sequence labeling is to find the best way to assign a categorical label to each token of a sequence. Formally, for a generatic input sentence $\mathbf{x} = \{w_1, \cdots, w_n\}$, $\mathbf{y} = \{y_1, \cdots, y_n\}$ represents a generic sequence of labels for $\mathbf{x}$.

Two classic tasks in sequence labeling, which are also the tasks considered in this thesis, are part-of-speech (POS) tagging and named entity recognition (NER). These tasks are one of the first stages in deep language understanding, whose importance has been well recognized in the natural language processing (NLP) community. As mentioned in Chapter 1, POS tagging is to assign to each word in a sentence its part-of-speech tag to represent its syntac function, while NER is to recognize the spans of named entities in the given sentence. For POS tagging, it is straight-froward to formulate it as a sequence labeling problem. For NER, a typical way to set this

up as sequence labeling is to use the *BIO tagging schema* (as shown in Figure 2.2). Each word is labeled *B* (beginning) if it is the first word in a named entity, *I* (inside) if it is a subsequent word in a named entity, or *O* (outside) otherwise.

As to the techniques applied in sequence labeling, there are four main streams: 1) Rule-based approaches, which rely on hand-crafted rules; 2) Unsupervised learning approaches which rely on unsupervised algorithms without hand-labeled training instances; 3) Feature-based supervised learning approaches, which rely on supervised learning algorithms with carefully designed features; 4) Deep learning approaches, which automatically discover feature representations from raw input in an end-to-end manner. The first three categories of approaches, which are also called traditional or classical approaches, rely on expertise on specific tasks (either in the form of manual feature engineering or hand-crafted rules). Deep learning approaches, on the other hand, is able to automatically extract salient features for downstream tasks but require large amount of training instances. In this section, we briefly describe the traditional approaches for sequence labeling, leaving the discussion of deep learning approaches to Section 2.4.3.

### 2.2.1 Rule-based Approaches

Rule-based sequence labeling systems rely on hand-crafted rules that are designed based on syntactic-lexical patterns (Brill, 1992; Sekine and Nobata, 2004; Etzioni et al., 2005) and domain-specific gazetteers (Zhang and Elhadad, 2013). Kim and Woodland (2000) proposed a system that generates rules automatically based on Brill's par-tof-speech tagger (Brill, 1992). In biomedical domain, Hanisch et al. (2005) proposed ProMiner, which leverages a pre-processed synonym dictionary to identify protein mentions and potential gene in biomedical text. Quimbaya et al. (2016) developed a dictionary-based approach for NER in electronic health records. Some other well-known rule-based NER systems, which are mainly based on hand-crafted syntactic and semantic rules to recognize entities, include LaSIE-II (Humphreys et al., 1998), LTG (Mikheev et al., 1999), and NewOwl (Krupka and IsoQuest, 2005). Rule-based systems work very well when lexicon is exhaustive. Due to domain-specific rules and incomplete dictionaries, high precision and low recall are often observed from such systems, and the systems cannot be transferred to other domains (Li et al., 2020).

## 2.2.2 Unsupervised Learning Approaches

Unsupervised sequence labeling is a classic problem in unsupervised learning that has been explored with various approaches. For unsupervised POS tagging, Haghighi and Klein (2006) assumed a set of prototypical words for each tag and report high accuracy. Johnson (2007) investigated why the hidden Markov models (HMMs) estimated by Expectation-Maximization (EM) produce poor results as POS taggers. Berg-Kirkpatrick et al. (2010) proposed an HMM in which probabilties are given by log-linear models. Stratos et al. (2016) tackled unsupervised POS tagging with HMMs by imposing an assumption that each hidden state is associated with an observation state that can appear under no other state. Some other methods leverage additional results. For example, Das and Petrov (2011) and Täckström et al. (2013) utilized parallel data to project POS tags from a source language. Li et al. (2012) used tag dictionaries from Wiktionary.

For NER, a typical approach of unsupervised learning is clustering (Nadeau and Sekine, 2007), where clustering-based systems extract labels from the clustered groups based on context similarity. The key idea is that lexical resources, lexical patterns, and statistics computed on a large corpus can be used to infer mentions of named entities (Li et al., 2020). Collins and Singer (1999) used unlabeled data to reduce the requirements for supervision and presented two unsupervised algorithms for named entity classification. Nadeau et al. (2006) proposed an unsupervised system for gazetteer building and named entity ambiguity resolution. Zhang and Elhadad (2013) proposed an unsupervised approach to extract named entities from biomedical text. Instead of supervision, their model resorts to terminologies, corpus statistics and shallow syntactic knowledge.

## 2.2.3 Feature-based Supervised Learning Approaches

Applying supervised learning, feature engineering is critical in supervised sequence labeling systems. Given annotated data samples, features are carefully designed to represent each training example. Machine learning algorithms are then utilized to learn a model to recognize similar patterns from unseen data. Feature vector representation is an abstraction over text where a word is represented by one or many Boolean, numeric, or nominal values (Nadeau and Sekine, 2007). Word-level features, dictionary lookup features, and document and corpus features have been

widely used in various supervised sequence labeling systems. More feature designs are discussed in (Nadeau and Sekine, 2007; Manning, 2011; Campos et al., 2012; Sharnagat, 2014).

Based on these features, many machine learning algorithms have been applied in supervised sequence labeling models. Crammer and Singer (2001) proposed the multiclass SVM method that casts to a multi-class classification task. HMMs (Rabiner, 1989) are a traditional statistical tool for modeling sequences and have been widely applied to linguistic sequence labeling tasks, such as POS tagging and NER (Bikel et al., 1997, 1999). Lafferty et al. (2001) introduced conditional random fields (CRF) to sequence labeling. Since then, CRF-based models have been widely applied to sequence labeling tasks in various domains, including biomedical text(Settles, 2004; Liu et al., 2019b), chemical text (Rocktäschel et al., 2012) and tweets (Ritter et al., 2011; Liu et al., 2011), and spawned many variants (McCallum and Li, 2003; Krishnan and Manning, 2006).

## 2.3    Dependency Parsing

Syntactic dependency representations have a long history in descriptive and theoretical linguistics and many formal models have been advanced, most notably Word Grammar (Hudson, 1984), Functional Generative Description (Sgall et al., 1986), Meaning-Text Theory (Melćuk et al., 1988), and Constraint Dependency Grammar (Maruyama, 1990). Syntactic dependency graphs have gained a wide interest in the computational linguistics community and have been successfully employed for a wide range of NLP applications such as entity coreference resolution (Ng, 2010; Durrett and Klein, 2013; Ma et al., 2016), sentiment analysis(Tai et al., 2015), machine translation (Bastings et al., 2017), and information extraction (Nguyen et al., 2009; Angeli et al., 2015; Peng et al., 2017).

Dependency trees represent syntactic relationships between words in the sentences through labeled directed edges between head words and their dependents (modifiers). The task of dependency parsing is to automatically analyze the dependency structure for a given sentence. Figure 2.3 shows a dependency tree for the sentence *I saw a girl with a telescope.*, with the symbol $ as its root. Dependency trees are often typed with labels for each edge to represent additional syntactic information, such as *sbj* and *dobj* for verb-subject and verb-object head-modifier interactions,

Figure 2.3: An example of dependency tree structure with dependency labels on edges.

respectively. Sometimes, however, the dependency labels are omitted. Dependency trees are defined as labeled or unlabeled according to whether the dependency labels are included or dropped. In the remainder of this thesis, we will focus on labeled dependency parsing, i.e. jointly predicting the dependency tree structures and the dependency labels in a single parsing model.

By considering the item of crossing dependencies, dependency trees fall into two categories — projective and non-projective dependency trees. An equivalent and more convenient formulation of the projectivity constrain is that if a dependency tree can be written with all words in a predefined linear order and all edges (including the edge for root) drawn on the plane without crossing edges. The example in Figure 2.3 belongs to the class of projective dependency trees. Previous studies illustrate that projective dependency trees are sufficient to analyze most English sentences, due to English's rigid word order (Yamada and Matsumoto, 2003; McDonald et al., 2005a). However, for languages with flexible word orders, non-projective dependency tree is preferable. In this thesis, we consider non-projective dependency parsing, where we develop dependency parsers that is not restricted by the projective constraint. More detailed discussion on dependency parsing can be found in (McDonald, 2006).

In this thesis, we focus on data-driven dependency parsing supported by supervised learning algorithms. According to the CoNLL-X shared tasks on dependency parsing (Buchholz and Marsi, 2006; Nivre et al., 2007), there are two dominant approaches to dependency parsing (Buchholz and Marsi, 2006; Nivre et al., 2007): local and greedy *transition-based* algorithms (Yamada and Matsumoto, 2003; Nivre and Scholz, 2004; Zhang and Nivre, 2011; Chen and Manning, 2014), and the globally optimized *graph-based* algorithms (Eisner, 1996; McDonald et al., 2005a,b; Koo and Collins, 2010). However, as discussed in McDonald and Nivre (2011), there is no a priori

reason why a graph-based parameterization should require global learning and inference,and a transition-based parameterization would necessitate local learning and greedy inference. In the following of this section, we will briefly review these two parsing models.

## 2.3.1 Graph-based Dependency Parsing

Formally, for a given sentence $\mathbf{x} = \{w_1, \ldots, w_n\}$, the complete dependency graph $G(\mathbf{x}) =< V(\mathbf{x}), E(\mathbf{x}) >$ represent words and their relationship to syntactic modifiers using directed edges. $V(\mathbf{x})$ and $E(\mathbf{x})$ are the sets of vertex and edges of the graph $G(\mathbf{x})$. A valid dependency tree $\mathbf{y}$ of $\mathbf{x}$ can be seen as a subgraph of $G(\mathbf{x})$ that build up a tree structure. We write $(w_i \rightarrow w_j) \in \mathbf{y}$ if there is a dependency in $\mathbf{y}$ from word $w_i$ to word $w_j$. Graph-based dependency parsers learn scoring functions for parse trees and perform exhaustive search over all possible trees for a sentence to find the globally highest scoring tree.

This category of models parameterize over dependency subgraphs and learns these parameters to globally score correct graphs above incorrect ones. Inference is also global,in that systems attempt to find the highest scoring graph among the set of all graphs. We call such systems *graph-based* parsing models to reflect the fact that parameterization is over the dependency graph. Eisner (1996) gave a generative model with a cubic parsing algorithm based on a graph factorization that inspired the core algorithms for graph-based dependency parsing. Following this pioneering work, several graph-based dependency parsing models are proposed, together with associated learning and decoding algorithms (McDonald et al., 2005a,b; McDonald and Pereira, 2006; Koo et al., 2007; Smith and Smith, 2007).

**Edge-Factored Parsing Model.** A common method is to factorize the score of a dependency tree as the sum of the scores of all edges in the tree:

$$\psi(\mathbf{x}, \mathbf{y}; \theta) = \sum_{(w_h, w_m) \in \mathbf{y}} \psi(w_h, w_m; \theta)$$

where $\theta$ is the parameter and $\psi(\mathbf{x}, \mathbf{y}; \theta)$ is the score function of the parse tree $\mathbf{y}$, which is factorized as the sum of the scores of each edge $\psi(w_h, w_m; \theta)$.

**Maximum Spanning Tree Decoding.** The decoding problem of this parsing model can be formulated as:

$$
\begin{aligned}
\mathbf{y}^* &= \operatorname*{argmax}_{\mathbf{y} \in T(\mathbf{x})} \psi(\mathbf{y}|\mathbf{x}; \theta) \\
&= \operatorname*{argmax}_{\mathbf{y} \in T(\mathbf{x})} \sum_{(w_h, w_m) \in \mathbf{y}} \psi(w_h, w_m; \theta)
\end{aligned}
$$

which can be solved by using the Maximum Spanning Tree (MST) algorithm described in McDonald et al. (2005b).

**Higher-order Factorizations.** A common strategy to improve the edge-factorization is to utilize high-order factorization:

$$
\psi(\mathbf{x}, \mathbf{y}; \theta) = \sum_{p \in \mathbf{y}} \psi(p; \theta)
$$

where $p$ is a part of the dependency tree $\mathbf{y}$.



McDonald and Pereira (2006) proposed to factorize each tree into second-order *sibling* parts — parts of dependencis consists of a triple of indices $(h, m, s)$ where $(h, m)$ and $(h, s)$ are dependencies, and where $s$ and $m$ are successive modifiers to the same side of $h$. Carreras (2007) first proposed to use a second-order *grandchild* part — pairs of dependencies connected head-to-tail. Koo and Collins (2010) introduced a third-order *grand-sibling* part — combinations of sibling parts and grandchild parts. Further, Ma and Zhao (2012a) combined grand-sibling and tri-sibling parts to propose the fourth-order *grand-tri-sibling* part.

## 2.3.2 Transition-based Dependency Parsing

The transition-based parsing systems parameterize over transitions from one state to another in an abstract state-machine. Formally, a transition system for dependency parsing defines

- **Parser configurations**: a set $C$ contains a (partially built) dependency graph $G$.

- **Transitions:** a set $T$ of transitions, each of which is a partial function $t : C \to C$.

- **Initial and terminal configurations:** for every sentence **x**, a unique initial configuration $c_x$ and a set of terminal configurations $C_x$.

Transition-based parsers learn scoring functions $s : C \times T \to \mathcal{R}$ that represent the likelihood of taking transition $t$ out of configuration $c$ in a transition sequence leading to the optimal dependency graph for the given sentence. Parameters in these score functions are typically learned using standard classification techniques that learn to predict one transition from a set of permissible transitions given a state history. We call such systems *transition-based* parsing models to reflect the fact that parameterization is over possible state transitions.

Inference is local, in that systems start in a fixed initial configuration $c_x$ and greedily construct the graph by taking the highest scoring transitions at each state entered until a termination configuration $c_m \in C_x$ is met:

$$t^* = \operatorname*{argmax}_{t \in T} s(c, t) \tag{2.4}$$

This can be seen as a greedy search for the optimal dependency graph,based on a sequence of locally optimal decisions in terms of the transition system.

Nivre (2003) introduced the arc-eager transition system that can derive any projective dependency tree for an input sentence. When coupled with the greedy deterministic parsing strategy (Yamada and Matsumoto, 2003; Nivre and Scholz, 2004), the system guarantees termination after at most $2n$ transitions (linear w.r.t the sentence length $n$). Nivre et al. (2004) extended this system to labeled dependency graphs, and Nivre and Nilsson (2005) showed how to lift the restriction of projective dependency graphs by using the pseudo-projective parsing technique. Transition systems that derive non-projective trees directly have been explored by Attardi (2006); Nivre (2007, 2009).

## 2.4   Feature Engineering vs. Representation Learning

To make task-oriented predictions, the development of machine learning systems heavily rely on extracting abstract informative features from real-world data that are represented in raw digital

a)

| Basic Uni-gram Features |
| --- |
| $x_i$-word, $x_i$-pos |
| $x_i$-word |
| $x_i$-pos |
| $x_j$-word, $x_j$-pos |
| $x_j$-word |
| $x_j$-pos |

b)

| Basic Bi-gram Features |
| --- |
| $x_i$-word, $x_i$-pos, $x_j$-word, $x_j$-pos |
| $x_i$-pos, $x_j$-word, $x_j$-pos |
| $x_i$-word, $x_j$-word, $x_j$-pos |
| $x_i$-word, $x_i$-pos, $x_j$-pos |
| $x_i$-word, $x_i$-pos, $x_j$-word |
| $x_i$-word, $x_j$-word |
| $x_i$-pos, $x_j$-pos |

c)

| In Between POS Features |
| --- |
| $x_i$-pos, b-pos, $x_j$-pos |
| **Surrounding Word POS Features** |
| $x_i$-pos, $x_i$-pos+1, $x_j$-pos-1, $x_j$-pos |
| $x_i$-pos-1, $x_i$-pos, $x_j$-pos-1, $x_j$-pos |
| $x_i$-pos, $x_i$-pos+1, $x_j$-pos, $x_j$-pos+1 |
| $x_i$-pos-1, $x_i$-pos, $x_j$-pos, $x_j$-pos+1 |

d)

| Second-order Features |
| --- |
| $x_i$-pos, $x_k$-pos, $x_j$-pos |
| $x_k$-pos, $x_j$-pos |
| $x_k$-word, $x_j$-word |
| $x_k$-word, $x_j$-pos |
| $x_k$-pos, $x_j$-word |

Figure 2.4: Features used by the MST-Parser (McDonald and Pereira, 2006) for a single dependency edge.

formats. For example, syntactic, semantic and contextual information are essentially important for a wide range of NLP tasks such as information extraction (IE). Text documents, however, are stored by individual tokens of words or characters, upon which these abstract information is hard to represent. Similar scenarios happen in image data — digital images are made up of pixels, from which it is difficult to extract abstract features such as edge or shape that are crucial for tasks, for instance, image classification.

## 2.4.1 Feature Engineering

Feature engineering is critical in traditional supervised machine learning systems. In the literature of linguistic structured prediction, most traditional high performance structured prediction models extract instructive features using the hand-crafted feature-engineering process, which mainly relies on task-specific expertise and heuristically designed hand-crafted features, in an iterative feature-selection process. Word-level features (e.g., case, morphology, and part-of-speech tag) (Zhou and Su, 2002; Settles, 2004; Liao and Veeramachaneni, 2009), list lookup features (e.g., Wikipedia gazetteer and DBpedia gazetteer) (Mikheev, 1999; Torisawa et al., 2007; Hoffart et al., 2011), and document and corpus features (e.g., local syntax and multiple occurrences) (Ravin and Wacholder, 1997; Krishnan and Manning, 2006) have been widely used in various structured prediction systems. For example, English POS taggers benefit from carefully designed word spelling features; orthographic features and external resources such as gazetteers are widely used in NER. Figure 2.4 lists the concrete features used in the MST-Parser (McDonald and Pereira, 2006) to represent a single dependency edge $w_i \rightarrow w_j$. Table (a) and (b) show the basic set of features, which are over head-modifier pairs in the tree. These features provide back-off from

22

very specific features over words and part-of-speech (POS) tags to less sparse features over just POS tags. These features are added for both the entire words as well as the 5-gram prefix if the word is longer than 5 characters (McDonald et al., 2005a). In addition to the basic set of features, McDonald and Pereira (2006) added two more types of features, listed in Table (c). The first new feature class recognizes word types that occur between the head and modifier words in an attachment decision and the second class of additional features represents the local context of the attachment, that is, the words before and after the head-modifier pair. All these features are carefully designed with linguistic knowledge and selected using the iterative selection process.

Despite the successes of feature-based approaches on linguistic structured prediction, there are two problems with that brute-force methodology: 1) the combinatorial nature of empirical feature selection process makes it expensive to hand-craft features; 2) The development of these features is commonly task-, domain-, or even language-specific, preventing it from adapting to new tasks or domains.

## 2.4.2 End-to-end Representation Learning

In the past few years, deep neural networks, together with end-to-end learning paradigm, have been applied to a wide range of NLP tasks with great successes. Unlike traditional linear models that rely on manual feature engineering, deep neural models are able to learn complex and intricate features from data via non-linear activation functions. The feature representation $\phi(w_i)$ is no longer binary vectors, but continuous vectors output from neural networks. Furthermore, deep neural models can be trained in an end-to-end paradigm, by back-propagation. This saves significant effort on designing task-specific features.

Formally, we use $\phi_{\theta'}(\mathbf{x})$ to denote the neural representation of $\mathbf{x}$ output from a neural network $\mathbf{M}_{\theta'}$, where $\theta'$ are the parameters of the neural network $\mathbf{M}$. So in deep learning models, the model parameters are the union of the parameters of the machine learning model and the neural network: $\theta \cup \theta'$. During model training the two sets of parameters are updated simultaneously to optimize the loss function in an end-to-end learning paradigm:

$$\min_{\theta, \theta' \in \Theta} \frac{1}{N} \sum_{i=1}^{N} - \log P_{\theta}(\mathbf{y}_i | \phi_{\theta'}(\mathbf{x}_i))$$

In the rest of this thesis, we use $\theta$ to denote the set of all the parameters of a deep learning model, when there is no ambiguity.

The end-to-end training paradigm significantly simplifies the hand-crafted feature engineering process in traditional feature-based machine learning systems, while giving the neural models flexibility to be optimized towards the ultimate tasks.

### 2.4.3    A Brief History of Neural Representation Learning on Linguistic Structured Prediction

**Distributed representations for input.**    Distributed representation represents words in low dimensional real-valued dense vectors where each dimension represents a latent feature. Automatically learned from text, distributed representation captures semantic and syntactic properties of word. Next, we first review neural structured predictions models with distributed representations.

Henderson (2004) first attempted to use neural networks to predict parse decisions in a constituency parser. Titov and Henderson (2007a) developed a generative dependency parser with latent variable models. Titov and Henderson (2007b) applied Incremental Sigmoid Belief Networks to constituency parsing and then Garg and Henderson (2011) extended this work to transition-based dependency parsers using a Temporal Restricted Boltzman Machine. More recently, Chen and Manning (2014) proposed to use feed-forward neural networks to model the transition states. For sequence labeling, Collobert et al. (2011) proposed a simple but effective feed-forward neutral network that independently classifies labels for each word by using contexts within a window with fixed size. This model achieved impressive performance on sequence labeling tasks, including POS tagging and NER.

**RNN-based context encoding architecture.**    Recurrent neural networks (RNNs), together with its variants such as gated recurrent unit (GRU) and long-short term memory (LSTM), have demonstrated remarkable achievements in modeling sequential data. Graves and Schmidhuber (2005) first attempted to apply LSTM to phoneme classification and Graves (2008) applied RNNs to sequence labeling. Recently, several RNN-based neural network models have been proposed to solve sequence labeling tasks like speech recognition (Graves et al., 2013), POS tagging (Santos

24

and Zadrozny, 2014; Huang et al., 2015; Labeau et al., 2015) and NER (dos Santos et al., 2015; Hu et al., 2016; Peng and Dredze, 2016; Lample et al., 2016; Chiu and Nichols, 2016), achieving competitive performance against traditional models.

For graph-based dependency parsing, Kiperwasser and Goldberg (2016) and Wang and Chang (2016) replaced the linear scoring function of each arc in traditional models with neural networks and used a margin-based objective McDonald et al. (2005a) for model training. Kiperwasser and Goldberg (2016)proposed a graph-based dependency parser which uses bidirectional LSTM for word-level representations. Wang and Chang (2016) used a similar model with a way to learn sentence segment embedding based on an extra forward LSTM network. Both of these two parsers trained the parsing models by optimizing margin-based objectives. Zhang et al. (2016) and Dozat and Manning (2017) formalized dependency parsing as independently selecting the head of each word with cross-entropy objective, without the guarantee of a non-projective structure.

For transition-based dependency parsing, neural continuous states have been explored, in which the transition state is embedded as a neural continuous vector. Dyer et al. (2015) introduced transtion-based parser using Stack LSTMs whose continuous-state embeddings were constructed using LSTM recurrent neural networks which are capable of learning representations of the parser state that are sensitive to the complete contents of the parser's state. Ballesteros et al. (2015) improved the Stack-LSTM parser by replacing word representations with representations constructed from the orthographic representations of the words, and Ballesteros et al. (2016) adapted the Stack-LSTM parser to support training-with-exploration procedure using dynamic oracles. Weiss et al. (2015a) presented structured perceptron training for neural network transition-based dependency parsing. Andor et al. (2016) proposed a globally normalized transition model to replace the locally normalized classifier.

**Pre-trained contextualized word representations.**    Recently, deep contextualized representations pre-trained on large scale corpus with language-model-based objectives have been empirically verified to be helpful in numerous structured prediction tasks. Peters et al. (2017) proposed TagLM, a language model augmented sequence tagger. This tagger considers both pre-trained word embeddings and bidirectional language model embeddings for every token in the input sequence for sequence labeling task. Following this work, Peters et al. (2018) proposed a new

type of deep contextualized word representation, named ELMo, which is capable of capturing both complex characteristics of word usage (e.g., syntax and semantics), and usage variations across linguistic contexts (e.g., polysemy). Based-on transformer architecture (Vaswani et al., 2017), Devlin et al. (2019) proposed bidirectional encoder representations (BERT). Several studies (Clark et al., 2018; Luo et al., 2019; Liu et al., 2019c) have achieved promising performance via leveraging deep contextualized representations.

## 2.5   BLSTM-CNNs

In this section, we describe our BLSTM-CNNs neural encoding architecture which consistently obtains impressive performance across different structured prediction tasks. Specificiially, we first use convolutional neural networks (CNNs) LeCun et al. (1989) to encode character-level information of a word into its character-level representation. Then we combine character- and word-level representations and feed them into bi-directional LSTM (BLSTM) to model context information of each word. It is a truly end-to-end model requiring no task-specific resources, feature engineering, or data pre-processing beyond pre-trained word embeddings on unlabeled corpora. Thus, our model can be easily applied to a wide range of structured prediction tasks on different languages and domains (see Chapter 3 and 4).

The main contribution of BLSTM encoder is that it is one of the first neural architecture for encoding sentences, which captures both word- and character-level information. Furthermore, it is truly end-to-end, entirely getting rid of feature engineering, and achieved state-of-the-art or comparable performance across various structured prediction tasks. It is an important step towards end-to-end representation learning for texts and inspires following works such as deep contextualized representations (Peters et al., 2018). In the following sections, we describe the components (layers) of our neural network architecture one-by-one from bottom to top.

### 2.5.1   CNN for Character-level Representation

Previous studies (Santos and Zadrozny, 2014; Chiu and Nichols, 2016) have shown that CNN is an effective approach to extract morphological information (like the prefix or suffix of a word)

Figure 2.5: The convolution neural network for extracting character-level representations of words. Dashed arrows indicate a dropout layer applied before character embeddings are input to CNN.

from characters of words and encode it into neural representations. Figure 2.5 shows the CNN we use to extract character-level representation of a given word.

The CNN is similar to the one in Chiu and Nichols (2016), except that we use only character embeddings as the inputs to CNN, without character type features. A dropout layer (Srivastava et al., 2014b) is applied before character embeddings are input to CNN.

### 2.5.2 Bi-directional LSTM

**LSTM Unit**

Recurrent neural networks (RNNs) are a powerful family of connectionist models that capture time dynamics via cycles in the graph, and are capable of dealing with variable-length sequence input. It uses a recurrent hidden state whose activation is dependent on that of the one immediate before. Though, in theory, RNNs are capable to capturing long-distance dependencies, in practice, they fail due to the gradient vanishing/exploding problems (Bengio et al., 1994; Pascanu et al.,

2012). In order to mitigate this weak point in conventional RNNs, specially designed activation functions have been introduced. LSTMs (Hochreiter and Schmidhuber, 1997) are variants of RNNs designed to cope with these gradient vanishing problems. Basically, a LSTM unit is composed of three multiplicative gates which control the proportions of information to forget and to pass on to the next time step. Figure 2.6 gives the basic structure of an LSTM unit.



Figure 2.6: Schematic of LSTM unit.

Formally, the formulas to update an LSTM unit at time $t$ are:

$$
\begin{aligned}
\mathbf{i}_t &= \sigma(\boldsymbol{W}_i\mathbf{h}_{t-1} + \boldsymbol{U}_i\mathbf{x}_t + \boldsymbol{b}_i) \\
\mathbf{f}_t &= \sigma(\boldsymbol{W}_f\mathbf{h}_{t-1} + \boldsymbol{U}_f\mathbf{x}_t + \boldsymbol{b}_f) \\
\tilde{\mathbf{c}}_t &= \tanh(\boldsymbol{W}_c\mathbf{h}_{t-1} + \boldsymbol{U}_c\mathbf{x}_t + \boldsymbol{b}_c) \\
\mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \\
\mathbf{o}_t &= \sigma(\boldsymbol{W}_o\mathbf{h}_{t-1} + \boldsymbol{U}_o\mathbf{x}_t + \boldsymbol{b}_o) \\
\mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t)
\end{aligned}
$$

where $\sigma$ is the element-wise sigmoid function and $\odot$ is the element-wise product. $\mathbf{x}_t$ is the input vector (e.g. word embedding) at time $t$, and $\mathbf{h}_t$ is the hidden state (also called output) vector storing all the useful information at (and before) time $t$. $\boldsymbol{U}_i, \boldsymbol{U}_f, \boldsymbol{U}_c, \boldsymbol{U}_o$ denote the weight matrices of different gates for input $\mathbf{x}_t$, and $\boldsymbol{W}_i, \boldsymbol{W}_f, \boldsymbol{W}_c, \boldsymbol{W}_o$ are the weight matrices for hidden state $\mathbf{h}_t$. $\boldsymbol{b}_i, \boldsymbol{b}_f, \boldsymbol{b}_c, \boldsymbol{b}_o$ denote the bias vectors.

**Convolution Layers**

**LSTM Layers**



Figure 2.7: The main architecture of our encoding neural network. The character representation for each word is computed by the CNN in Figure 2.5. Then the character representation vector is concatenated with the word embedding before feeding into the BLSTM network. Dashed arrows indicate dropout layers applied on both the input and output vectors of BLSTM.

LSTM uses input and output gates to control the flow of information through the cell. The input gate should be kept sufficiently active to allow the signals in. Same rule applies to the output gate. The forget gate is used to reset the cell's own state. In Gers et al. (2003), peephole connections are usually used to connect gates to the cell in tasks requiring precise timing and counting of the internal states. It should be noted that we do not include peephole connections (Gers et al., 2003) in the our LSTM formulation because the precise timing does not seem to be required.

**BLSTM**

Another weak point of conventional RNNs is their utilization of only previous context with no exploitation of future context. While for many linguistic structured prediction tasks it is beneficial to have access to both past (left) and future (right) contexts. However, the LSTM's hidden state $\mathbf{h}_t$ takes information only from past, knowing nothing about the future. An elegant solution whose effectiveness has been proven by previous work (Dyer et al., 2015) is bi-directional LSTM

(BLSTM). The basic idea is to present each sequence forwards and backwards to separate hidden states to capture both past and future information. Then the two hidden states are concatenated to form the final output.

### 2.5.3   BLSTM-CNNs Encoding Architecture

For each word, the character-level representation is computed by the CNN in Figure 2.5 with character embeddings as inputs. Then the character-level representation vector is concatenated with the word embedding vector to feed into the BLSTM network.  Figure 2.7 illustrates the architecture of our network in detail.

# Chapter 3

# Sequence Labeling via BLSTM-CNNs-CRF

This chapter describes the BLSTM-CNNs-CRF model for linguistic sequence labeling.

## 3.1 Neural CRF Model

Following the notations defined in Section 2.1, we use $\mathbf{x} = \{w_1, \cdots, w_n\}$ to represent a generic input sequence where $w_i$ is the input vector of the $i$th word. $\boldsymbol{y} = \{y_1, \cdots, y_n\}$ represents a generic sequence of labels for $\mathbf{x}$. $T(\mathbf{x})$ denotes the set of possible label sequences for $\mathbf{x}$.

A simple and straight-forward solution to model the conditional probability $P_\theta(\mathbf{y}|\mathbf{x})$ over all $\mathbf{y}$ is to assume that each token of label sequence is independent given the input sentence $\mathbf{x}$:

$$P_\theta(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{n} P_\theta(y_i|\mathbf{x})$$

For sequence labeling (or general structured prediction) tasks, however, it is beneficial to consider the correlations between labels in neighborhoods and jointly decode the best chain of labels for a given input sentence. For example, in POS tagging an adjective is more likely to be followed by a noun than a verb, and in NER with standard BIO2 annotation (Tjong Kim Sang and Veenstra, 1999) I-ORG cannot follow I-PER. Therefore, we model label sequence jointly using a conditional random field (CRF) (Lafferty et al., 2001), instead of decoding each label independently.

### 3.1.1  Conditional Random Fields (CRF)

The probabilistic model for sequence CRF (Lafferty et al., 2001) defines a family of conditional probability $p(\boldsymbol{y}|\mathbf{x}; \mathbf{W}, \mathbf{b})$ over all possible label sequences $\boldsymbol{y}$ given $\mathbf{x}$ with the following form:

$$P_\theta(\boldsymbol{y}|\mathbf{x}) = \frac{\prod_{i=1}^{n} \psi_i(y_{i-1}, y_i, \mathbf{x})}{\sum_{y' \in T(\mathbf{x})} \prod_{i=1}^{n} \psi_i(y'_{i-1}, y'_i, \mathbf{x})}$$

where $\psi_i(y', y, \mathbf{x}) = \exp(\mathbf{W}_{y',y}^T \phi(w_i) + \mathbf{b}_{y',y})$ are potential functions. $\phi(w_i)$ is the feature representation of word $w_i$. $\mathbf{W}_{y',y}^T$ and $\mathbf{b}_{y',y}$ are the weight vector and bias corresponding to label pair $(y', y)$, respectively. They are learnable parameters of the model, $\mathbf{W}_{y',y}^T, \mathbf{b}_{y',y} \in \theta$.

For a sequence CRF model (only interactions between two successive labels are considered), training and decoding can be solved efficiently by adopting the Viterbi algorithm (Forney, 1973; Lafferty et al., 2001).

### 3.1.2  BLSTM-CNNs-CRF

To construct our neural network model for sequence labeling, we feed the output vectors of BLSTM-CNNs into a CRF layer. Figure 3.1 illustrates the architecture of our network in detail.

For each word, the character-level representation is computed by the CNN in Figure 2.5 with character embeddings as inputs. Then the character-level representation vector is concatenated with the word embedding vector to feed into the BLSTM network. Finally, the output vectors of BLSTM ($\phi(w_i), i = 1, \ldots, n$) are fed to the CRF layer to jointly decode the best label sequence. As shown in Figure 3.1, dropout layers are applied on both the input and output vectors of BLSTM. Experimental results show that using dropout significantly improve the performance of our model (see Section 3.4.4 for details).

## 3.2  Network Training

In this section, we provide details about training the neural network. We implement the neural network using the Theano library (Bergstra et al., 2010). The computations for a single model are

Figure 3.1: The main architecture of our neural network. The character representation for each word is computed by the CNN in Figure 2.5. Then the character representation vector is concatenated with the word embedding before feeding into the BLSTM network. Dashed arrows indicate dropout layers applied on both the input and output vectors of BLSTM.

run on a GeForce GTX TITAN X GPU. Using the settings discussed in this section, the model training requires about 12 hours for POS tagging and 8 hours for NER.

### 3.2.1 Parameter Initialization

**Word Embeddings.** We use Stanford's publicly available GloVe 100-dimensional embeddings[1] trained on 6 billion words from Wikipedia and web text (Pennington et al., 2014)

We also run experiments on two other sets of published embeddings, namely Senna 50-dimensional embeddings[2] trained on Wikipedia and Reuters RCV-1 corpus (Collobert et al.,

[1] http://nlp.stanford.edu/projects/glove/
[2] http://ronan.collobert.com/senna/

2011), and Google's Word2Vec 300-dimensional embeddings[3] trained on 100 billion words from Google News (Mikolov et al., 2013). To test the effectiveness of pretrained word embeddings, we experimented with randomly initialized embeddings with 100 dimensions, where embeddings are uniformly sampled from range $[-\sqrt{\frac{3}{dim}}, +\sqrt{\frac{3}{dim}}]$ where $dim$ is the dimension of embeddings (He et al., 2015). The performance of different word embeddings is discussed in Section 3.4.3.

**Character Embeddings.** Character embeddings are initialized with uniform samples from $[-\sqrt{\frac{3}{dim}}, +\sqrt{\frac{3}{dim}}]$, where we set $dim = 30$.

**Weight Matrices and Bias Vectors.** Matrix parameters are randomly initialized with uniform samples from $[-\sqrt{\frac{6}{r+c}}, +\sqrt{\frac{6}{r+c}}]$, where $r$ and $c$ are the number of of rows and columns in the structure (Glorot and Bengio, 2010). Bias vectors are initialized to zero, except the bias $\mathbf{b}_f$ for the forget gate in LSTM , which is initialized to 1.0 (Jozefowicz et al., 2015).

### 3.2.2 Optimization Algorithm

Parameter optimization is performed with mini-batch stochastic gradient descent (SGD) with batch size 10 and momentum 0.9. We choose an initial learning rate of $\eta_0$ ($\eta_0 = 0.01$ for POS tagging, and $0.015$ for NER, see Section 3.2.3.), and the learning rate is updated on each epoch of training as $\eta_t = \eta_0/(1 + \rho t)$, with decay rate $\rho = 0.05$ and $t$ is the number of epoch completed. To reduce the effects of "gradient exploding", we use a gradient clipping of $5.0$ (Pascanu et al., 2012). We explored other more sophisticated optimization algorithms such as AdaDelta (Zeiler, 2012), Adam (Kingma and Ba, 2014) or RMSProp (Dauphin et al., 2015), but none of them meaningfully improve upon SGD with momentum and gradient clipping in our preliminary experiments.

**Early Stopping.** We use early stopping (Giles, 2001; Graves et al., 2013) based on performance on validation sets. The "best" parameters appear at around 50 epochs, according to our experiments.

[3]https://code.google.com/archive/p/word2vec/

| Layer | Hyper-parameter | POS | NER |
|---|---|---|---|
| CNN | window size | 3 | 3 |
| | number of filters | 30 | 30 |
| LSTM | state size | 200 | 200 |
| | initial state | 0.0 | 0.0 |
| | peepholes | no | no |
| Dropout | dropout rate | 0.5 | 0.5 |
| | batch size | 10 | 10 |
| | initial learning rate | 0.01 | 0.015 |
| | decay rate | 0.05 | 0.05 |
| | gradient clipping | 5.0 | 5.0 |

Table 3.1: Hyper-parameters for all experiments.

**Fine Tuning.**   For each of the embeddings, we fine-tune initial embeddings, modifying them during gradient updates of the neural network model by back-propagating gradients. The effectiveness of this method has been previously explored in sequential and structured prediction problems (Collobert et al., 2011; Peng and Dredze, 2015).

**Dropout Training.**   To mitigate overfitting, we apply the dropout method (Srivastava et al., 2014b) to regularize our model. As shown in Figure 2.5 and 3.1, we apply dropout on character embeddings before inputting to CNN, and on both the input and output vectors of BLSTM. We fix dropout rate at $0.5$ for all dropout layers through all the experiments. We obtain significant improvements on model performance after using dropout (see Section 3.4.4).

### 3.2.3   Tuning Hyper-Parameters

Table 3.1 summarizes the chosen hyper-parameters for all experiments. We tune the hyper-parameters on the development sets by random search. Due to time constrains it is infeasible to do a random search across the full hyper-parameter space. Thus, for the tasks of POS tagging and NER we try to share as many hyper-parameters as possible. Note that the final hyper-parameters

| Dataset | | WSJ | CoNLL2003 |
|---------|---------|---------|-----------|
| Train | SENT | 38,219 | 14,987 |
| | TOKEN | 912,344 | 204,567 |
| Dev | SENT | 5,527 | 3,466 |
| | TOKEN | 131,768 | 51,578 |
| Test | SENT | 5,462 | 3,684 |
| | TOKEN | 129,654 | 46,666 |

Table 3.2: Corpora statistics. SENT and TOKEN refer to the number of sentences and tokens.

for these two tasks are almost the same, except the initial learning rate. We set the state size of LSTM to 200. Tuning this parameter did not significantly impact the performance of our model. For CNN, we use 30 filters with window length 3.

## 3.3 Experiment Setup

### 3.3.1 Data Sets

As mentioned before, we evaluate our neural network model on two sequence labeling tasks: POS tagging and NER. The corpora statistics are shown in Table 3.2. We did not perform any pre-processing for data sets, leaving our system truly end-to-end.

**POS Tagging.** For English POS tagging, we use the Wall Street Journal (WSJ) portion of Penn Treebank (PTB) (Marcus et al., 1993), which contains 45 different POS tags. In order to compare with previous work, we adopt the standard splits — section 0–18 as training data, section 19–21 as development data and section 22–24 as test data (Manning, 2011; Søgaard, 2011).

**NER.** For NER, We perform experiments on the English data from CoNLL 2003 shared task (Tjong Kim Sang and De Meulder, 2003). This data set contains four different types of named entities: *PERSON, LOCATION, ORGANIZATION*, and *MISC*. We use the BIOES tagging scheme instead of standard BIO2, as previous studies have reported meaningful improvement with this scheme (Ratinov and Roth, 2009; Dai et al., 2015; Lample et al., 2016).

36

|  | POS | | NER | | | | | |
|---|---|---|---|---|---|---|---|---|
|  | Dev | Test | Dev | | | Test | | |
| **Model** | Acc. | Acc. | Prec. | Recall | F1 | Prec. | Recall | F1 |
| BRNN | 96.56 | 96.76 | 92.04 | 89.13 | 90.56 | 87.05 | 83.88 | 85.44 |
| BLSTM | 96.88 | 96.93 | 92.31 | 90.85 | 91.57 | 87.77 | 86.23 | 87.00 |
| BLSTM-CNN | 97.34 | 97.33 | 92.52 | 93.64 | 93.07 | 88.53 | 90.21 | 89.36 |
| BRNN-CNN-CRF | 97.46 | 97.55 | 94.85 | 94.63 | 94.74 | 91.35 | 91.06 | 91.21 |

Table 3.3: Performance of our model on both the development and test sets of the two tasks, together with three baseline systems.

## 3.4 Experimental Results

### 3.4.1 Main Results

We first run experiments to dissect the effectiveness of each component (layer) of our neural network architecture by ablation studies. We compare the performance with three baseline systems — BRNN, the bi-direction RNN; BLSTM, the bi-direction LSTM, and BLSTM-CNNs, the combination of BLSTM with CNN to model character-level information. All these models are run using Stanford's GloVe 100 dimensional word embeddings and the same hyper-parameters as shown in Table 3.1. According to the results shown in Table 3.3, BLSTM obtains better performance than BRNN on all evaluation metrics of both the two tasks. BLSTM-CNN models significantly [4] outperform the BLSTM model, showing that character-level representations are important for linguistic sequence labeling tasks. This is consistent with results reported by previous work (Santos and Zadrozny, 2014; Chiu and Nichols, 2016). Finally, by adding CRF layer for joint decoding we achieve significant improvements over BLSTM-CNN models for both POS tagging and NER on all metrics. This demonstrates that jointly decoding label sequences can significantly benefit the final performance of neural network models.

[4]We did not perform statistical significance test in this thesis. The significance of results is based on intuitive assessment of the magnitude of the difference.

| Model | Acc. | | Model | F1 |
|---|---|---|---|---|
| Giménez and Màrquez (2004) | 97.16 | | Chieu and Ng (2002) | 88.31 |
| Toutanova et al. (2003) | 97.27 | | Florian et al. (2003) | 88.76 |
| Manning (2011) | 97.28 | | Ando and Zhang (2005) | 89.31 |
| Collobert et al. (2011)‡ | 97.29 | | Collobert et al. (2011)‡ | 89.59 |
| Santos and Zadrozny (2014)‡ | 97.32 | | Huang et al. (2015)‡ | 90.10 |
| Shen et al. (2007) | 97.33 | | Chiu and Nichols (2016)‡ | 90.77 |
| Sun (2014) | 97.36 | | Ratinov and Roth (2009) | 90.80 |
| Søgaard (2011) | 97.50 | | Lin and Wu (2009) | 90.90 |
| **This paper** | **97.55** | | Passos et al. (2014) | 90.90 |
| | | | Lample et al. (2016)‡ | 90.94 |
| | | | Luo et al. (2015) | 91.20 |
| | | | **This paper** | **91.21** |

Table 3.4: Left: POS tagging accuracy of our model on test data from WSJ proportion of PTB, together with top-performance systems. Right: NER F1 score of our model on test data set from CoNLL-2003. For the purpose of comparison, we also list F1 scores of previous top-performance systems. The neural network based models are marked with ‡.

## 3.4.2 Comparison with Previous Work

**POS Tagging.** Table 3.4 (left) illustrates the results of our model for POS tagging, together with seven previous top-performance systems for comparison. Our model significantly outperform Senna (Collobert et al., 2011), which is a feed-forward neural network model using capitalization and discrete suffix features, and data pre-processing. Moreover, our model achieves 0.23% improvements on accuracy over the "CharWNN" (Santos and Zadrozny, 2014), which is a neural network model based on Senna and also uses CNNs to model character-level representations. This demonstrates the effectiveness of BLSTM for modeling sequential data and the importance of joint decoding with structured prediction model.

Comparing with traditional statistical models, our system achieves state-of-the-art accuracy,

obtaining 0.05% improvement over the previously best reported results by Søgaard (2011). It should be noted that Huang et al. (2015) also evaluated their BLSTM-CRF model for POS tagging on WSJ corpus. But they used a different splitting of the training/dev/test data sets. Thus, their results are not directly comparable with ours.

**NER.**    Table 3.4 (right) shows the F1 scores of previous models for NER on the test data set from CoNLL-2003 shared task. For the purpose of comparison, we list their results together with ours. Similar to the observations of POS tagging, our model achieves significant improvements over Senna and the other three neural models, namely the LSTM-CRF proposed by Huang et al. (2015), LSTM-CNNs proposed by Chiu and Nichols (2016), and the LSTM-CRF by Lample et al. (2016). Huang et al. (2015) utilized discrete spelling, POS and context features, Chiu and Nichols (2016) used character-type, capitalization, and lexicon features, and all the three model used some task-specific data pre-processing, while our model does not require any carefully designed features or data pre-processing. We have to point out that the result (90.77%) reported by Chiu and Nichols (2016) is incomparable with ours, because their final model was trained on the combination of the training and development data sets[5].

To our knowledge, the previous best F1 score (91.20)[6] reported on CoNLL 2003 data set is by the joint NER and entity linking model (Luo et al., 2015). This model used many hand-crafted features including stemming and spelling features, POS and chunks tags, WordNet clusters, Brown Clusters, as well as external knowledge bases such as Freebase and Wikipedia. Our end-to-end model slightly improves this model by 0.01%, yielding a state-of-the-art performance.

### 3.4.3   Word Embeddings

As mentioned in Section 3.2.1, in order to test the importance of pretrained word embeddings, we performed experiments with different sets of publicly published word embeddings, as well as a random sampling method, to initialize our model. Table 3.5 gives the performance of three different word embeddings, as well as the randomly sampled one. According to the results in

---

[5]We run experiments using the same setting and get 91.37% F1 score.

[6]Numbers are taken from the Table 3 of the original paper (Luo et al., 2015). While there is clearly inconsistency among the precision (91.5%), recall (91.4%) and F1 scores (91.2%), it is unclear in which way they are incorrect.

| Embedding | Dimension | POS | NER |
|---|---|---|---|
| Random | 100 | 97.13 | 80.76 |
| Senna | 50 | 97.44 | 90.28 |
| Word2Vec | 300 | 97.40 | 84.91 |
| GloVe | 100 | **97.55** | **91.21** |

Table 3.5: Results with different choices of word embeddings on the two tasks (accuracy for POS tagging and F1 for NER).

| | POS | | | NER | | |
|---|---|---|---|---|---|---|
| | **Train** | **Dev** | **Test** | **Train** | **Dev** | **Test** |
| No | 98.46 | 97.06 | 97.11 | 99.97 | 93.51 | 89.25 |
| Yes | 97.86 | 97.46 | 97.55 | 99.63 | 94.74 | 91.21 |

Table 3.6: Results with and without dropout on two tasks.

Table 3.5, models using pretrained word embeddings obtain a significant improvement as opposed to the ones using random embeddings. Comparing the two tasks, NER relies more heavily on pretrained embeddings than POS tagging. This is consistent with results reported by previous work (Collobert et al., 2011; Huang et al., 2015; Chiu and Nichols, 2016).

For different pretrained embeddings, Stanford's GloVe 100 dimensional embeddings achieve best results on both tasks, about 0.1% better on POS accuracy and 0.9% better on NER F1 score than the Senna 50 dimensional one. This is different from the results reported by Chiu and Nichols (2016), where Senna achieved slightly better performance on NER than other embeddings. Google's Word2Vec 300 dimensional embeddings obtain similar performance with Senna on POS tagging, still slightly behind GloVe. But for NER, the performance on Word2Vec is far behind GloVe and Senna. One possible reason that Word2Vec is not as good as the other two embeddings on NER is because of vocabulary mismatch — Word2Vec embeddings were trained in case-sensitive manner, excluding many common symbols such as punctuations and digits. Since we do not use any data pre-processing to deal with such common symbols or rare words, it might be an issue for using Word2Vec.

|      | POS      |          | NER   |       |
|------|----------|----------|-------|-------|
|      | Dev      | Test     | Dev   | Test  |
| IV   | 127,247  | 125,826  | 4,616 | 3,773 |
| OOTV | 2,960    | 2,412    | 1,087 | 1,597 |
| OOEV | 659      | 588      | 44    | 8     |
| OOBV | 902      | 828      | 195   | 270   |

Table 3.7: Statistics of the partition on each corpus. It lists the number of tokens of each subset for POS tagging and the number of entities for NER.

### 3.4.4 Effect of Dropout

Table 3.6 compares the results with and without dropout layers for each data set. All other hyper-parameters remain the same as in Table 3.1. We observe a essential improvement for both the two tasks. It demonstrates the effectiveness of dropout in reducing overfitting.

### 3.4.5 OOV Error Analysis

To better understand the behavior of our model, we perform error analysis on Out-of-Vocabulary words (OOV). Specifically, we partition each data set into four subsets — in-vocabulary words (IV), out-of-training-vocabulary words (OOTV), out-of-embedding-vocabulary words (OOEV) and out-of-both-vocabulary words (OOBV). A word is considered IV if it appears in both the training and embedding vocabulary, while OOBV if neither. OOTV words are the ones do not appear in training set but in embedding vocabulary, while OOEV are the ones do not appear in embedding vocabulary but in training set. For NER, an entity is considered as OOBV if there exists at lease one word not in training set and at least one word not in embedding vocabulary, and the other three subsets can be done in similar manner. Table 3.7 informs the statistics of the partition on each corpus. The embedding we used is Stanford's GloVe with dimension 100, the same as Section 3.4.1.

Table 3.8 illustrates the performance of our model on different subsets of words, together with the baseline LSTM-CNN model for comparison. The largest improvements appear on the OOBV

|              | POS | | | | | | | |
|              | Dev | | | | Test | | | |
|              | IV | OOTV | OOEV | OOBV | IV | OOTV | OOEV | OOBV |
|--------------|-------|-------|-------|-------|-------|-------|-------|-------|
| LSTM-CNN     | 97.57 | **93.75** | 90.29 | 80.27 | 97.55 | **93.45** | 90.14 | 80.07 |
| LSTM-CNN-CRF | **97.68** | 93.65 | **91.05** | **82.71** | **97.77** | 93.16 | **90.65** | **82.49** |
|              | NER | | | | | | | |
|              | Dev | | | | Test | | | |
|              | IV | OOTV | OOEV | OOBV | IV | OOTV | OOEV | OOBV |
| LSTM-CNN     | 94.83 | 87.28 | 96.55 | 82.90 | 90.07 | 89.45 | 100.00 | 78.44 |
| LSTM-CNN-CRF | **96.49** | **88.63** | **97.67** | **86.91** | **92.14** | **90.73** | 100.00 | **80.60** |

Table 3.8: Comparison of performance on different subsets of words.

subsets of both the two corpora. This demonstrates that by adding CRF for joint decoding, our model is more powerful on words that are out of both the training and embedding sets.

## 3.5 Discussions

In this chapter, we equipped the LSTM-CNNs encoding architecture with a CRF output layer for sequence labeling task. The evaluation results show that our truly end-to-end BLSTM-CNNs-CRF model achieved state-of-the-art performance on two linguistic sequence labeling tasks, comparing with previous state-of-the-art systems. The ablation studies analyzed the impact of each components of the model, indicating that the improvements come from both the encoding architecture and the CRF decoding mechanism.

# Chapter 4

# Neural Networks for Dependency Parsing

In this chapter, we introduce our NeuroMST parser for graph-based approach in Section 4.1 and leave the Stack-Pointer parser for transition-based approach in next chapter.

## 4.1 Neural Probabilistic Model for MST Parsing

### 4.1.1 Edge-Factored Probabilistic Model

The probabilistic model of NeuroMST parser defines a family of conditional probability $P_\theta(\mathbf{y}|\mathbf{x})$ over all valid parse trees $\mathbf{y}$ given a sentence $\mathbf{x}$, with a log-linear form:

$$P_\theta(\mathbf{y}|\mathbf{x}) = \frac{\exp\left(\sum_{(w_h, w_m) \in \mathbf{y}} \psi(w_h, w_m; \theta)\right)}{Z(\mathbf{x}; \theta)}$$

where $Z(\mathbf{x}; \theta)$ is the partition function.

$$Z(\mathbf{x}; \theta) = \sum_{\mathbf{y} \in T(\mathbf{x})} \exp\left(\sum_{(w_h, w_m) \in \mathbf{y}} \psi(w_h, w_m; \theta)\right)$$

**Bi-Linear Score Function.**   In our model, we adopt a bi-linear form score function:

$$\psi(w_h, w_m; \theta) = \phi(w_h)^T \mathbf{W} \phi(w_m) + \mathbf{U}^T \phi(w_h) + \mathbf{V}^T \phi(w_m) + \mathbf{b}$$

where $\{\mathbf{W}, \mathbf{U}, \mathbf{V}, \mathbf{b}\} \subset \theta$, $\phi(x_i)$ is the representation vector of $w_i$, $\mathbf{W}, \mathbf{U}, \mathbf{V}$ denote the weight matrix of the bi-linear term and the two weight vectors of the linear terms in $\psi$, and $\mathbf{b}$ denotes the bias vector.

As discussed in Dozat and Manning (2017), the bi-linear form of score function is related to the bi-linear attention mechanism (Luong et al., 2015). The bi-linear score function differs from the traditional score function proposed in Kiperwasser and Goldberg (2016) by adding the bi-linear term. A similar score function is proposed in Dozat and Manning (2017). The difference between their and our score function is that they only used the linear term for head words ($\mathbf{U}^T \phi(x_h)$) while use them for both heads and modifiers.

**Matrix-Tree Theorem.** In order to train the probabilistic parsing model, as discussed in Koo et al. (2007), we have to compute the *partition function* and the *marginals*, requiring summation over the set $T(\mathbf{x})$:

$$
\begin{aligned}
Z(\mathbf{x}; \theta) &= \sum_{\mathbf{y} \in T(\mathbf{x})} \exp \left( \sum_{(w_h, w_m) \in \mathbf{y}} \psi(w_h, w_m; \theta) \right) \\
\mu_{h,m}(\mathbf{x}; \theta) &= \sum_{\mathbf{y} \in T(\mathbf{x}):(w_h, w_m) \in \mathbf{y}} P(\mathbf{y}|\mathbf{x}; \theta)
\end{aligned}
$$

where $\mu_{h,m}(\mathbf{x}; \Theta)$ is the marginal for edge from $h$th word to $m$th word for $\mathbf{x}$.

Previous studies (Koo et al., 2007; Smith and Smith, 2007) have presented how a variant of Kirchhoff's Matrix-Tree Theorem (Tutte, 1984) can be used to evaluate the partition function and marginals efficiently. In this section, we briefly revisit this method.

For a sentence $\mathbf{x}$ with $n$ words, we denote $\mathbf{x} = \{w_0, w_1, \ldots, w_n\}$, where $w_0$ is the root-symbol. We define a complete graph $G$ on $n + 1$ nodes (including the root-symbol $x_0$), where each node corresponds to a word in $\mathbf{x}$ and each edge corresponds to a dependency arc between two words. Then, we assign non-negative weights to the edges of this complete graph with $n + 1$ nodes, yielding the weighted adjacency matrix $\mathbf{A}(\theta) \in \mathbb{R}^{n+1 \times n+1}$, for $h, m = 0, \ldots, n$:

$$
\mathbf{A}_{h,m}(\theta) = \exp \left( \psi(w_h, w_m; \theta) \right)
$$

Based on the adjacency matrix $\mathbf{A}(\theta)$, we have the Laplacian matrix:

$$
\mathbf{L}(\theta) = \mathbf{D}(\theta) - \mathbf{A}(\theta)
$$

where $\mathbf{D}(\theta)$ is the weighted degree matrix:

$$\mathbf{D}_{h,m}(\theta) = \begin{cases} \sum\limits_{h'=0}^{n} \mathbf{A}_{h',m}(\theta) & \text{if } h = m \\ 0 & \text{otherwise} \end{cases}$$

Then, according to Theorem 1 in Koo et al. (2007), the partition function is equal to the minor of $\mathbf{L}(\theta)$ w.r.t row 0 and column 0:

$$Z(\mathbf{x}; \theta) = \mathbf{L}^{(0,0)}(\theta)$$

where for a matrix $\mathbf{A}$, $\mathbf{A}^{(h,m)}$ denotes the *minor* of $\mathbf{A}$ w.r.t row $h$ and column $m$; i.e., the determinant of the submatrix formed by deleting the $h$th row and $m$th column.

The marginals can be computed by calculating the matrix inversion of the matrix corresponding to $\mathbf{L}^{(0,0)}(\theta)$. The time complexity of computing the partition function and marginals is $O(n^3)$.

**Labeled Parsing Model.**    Though it is originally designed for unlabeled parsing, our probabilistic parsing model is easily extended to include dependency labels.

In labeled dependency trees, each edge is represented by a tuple $(w_h, w_m, l)$, where $w_h$ and $w_m$ are the head word and modifier, respectively, and $l$ is the label of dependency type of this edge. Then we can extend the original model for labeled dependency parsing by extending the score function to include dependency labels:

$$\psi(w_h, w_m, l; \theta) = \phi(w_h)^T \mathbf{W}_l \phi(w_m) + \mathbf{U}_l^T \phi(w_h) + \mathbf{V}_l^T \phi(w_m) + \mathbf{b}_l$$

where $\mathbf{W}_l, \mathbf{U}_l, \mathbf{V}_l, \mathbf{b}_l$ are the weights and bias corresponding to dependency label $l$. Suppose that there are $L$ different dependency labels, it suffices to define the new adjacency matrix by assigning the weight of a edge with the sum of weights over different dependency labels:

$$\mathbf{A'}_{h,m}(\theta) = \sum_{l=1}^{L} \exp\left(\psi(w_h, w_m, l; \theta)\right)$$

The partition function and marginals over labeled dependency trees are obtained by operating on the new adjacency matrix $\mathbf{A'}(\theta)$. The time complexity becomes $O(n^3 + Ln^2)$. In practice, $L$ is probably large. For English, the number of edge labels in Stanford Basic Dependencies (De Marneffe et al., 2006) is 45, and the number in the treebank of CoNLL-2008 shared task (Surdeanu

et al., 2008) is 70, while the average length of sentences in English Penn Treebank (Marcus et al., 1993) is around 23. Thus, $L$ is not negligible to $n$.

It should be noticed that in our labeled model, for different dependency label $l$ we use the same vector representation $\phi(w_i)$ for each word $w_i$. The dependency labels are distinguished (only) by the parameters (weights and bias) corresponding to each of them. One advantage of this is that it significantly reduces the memory requirement comparing to the model in Dozat and Manning (2017) which distinguishes $\phi_l(w_i)$ for different label $l$.

## 4.1.2   Neural Representation Encoding

The encoder of our parsing model is based on the bi-directional LSTM-CNN architecture (BLSTM-CNNs) (Ma and Hovy, 2016) where CNNs encode character-level information of a word into its character-level representation and BLSTM models context information of each word. Formally, for each word, the CNN, with character embeddings as inputs, encodes the character-level representation. Then the character-level representation vector is concatenated with the word embedding vector to feed into the BLSTM network. To enrich word-level information, we also use POS embeddings. Figure 4.1 illustrates the architecture of our network in detail.

## 4.1.3   Neural Network Training

**Word Embeddings.**   For all the parsing models on different languages, we initialize word vectors with pretrained word embeddings. For Chinese, Dutch, English, German and Spanish, we use the structured-skipgram (Ling et al., 2015) embeddings, and for other languages we use the Polyglot (Al-Rfou et al., 2013) embeddings. The dimensions of embeddings are 100 for English, 50 for Chinese and 64 for other languages.

**Character Embeddings.**   Following Ma and Hovy (2016), character embeddings are initialized with uniform samples from $[-\sqrt{\frac{3}{dim}}, +\sqrt{\frac{3}{dim}}]$, where we set $dim = 50$.

**POS Embedding.**   Our model also includes POS embeddings. The same as character embeddings, POS embeddings are also 50-dimensional, initialized uniformly from $[-\sqrt{\frac{3}{dim}}, +\sqrt{\frac{3}{dim}}]$.

Figure 4.1: The main architecture of our parsing model. The character representation for each word is computed by the CNN in Figure 2.5. Then the character representation vector is concatenated with the word and pos embedding before feeding into the BLSTM network. Dashed arrows indicate dropout layers applied on the input, hidden and output vectors of BLSTM.

**Weights Matrices and Bias Vectors.** Matrix parameters are randomly initialized with uniform samples from $[-\sqrt{\frac{6}{r+c}}, +\sqrt{\frac{6}{r+c}}]$, where $r$ and $c$ are the number of of rows and columns in the structure (Glorot and Bengio, 2010). Bias vectors are initialized to zero, except the bias $\mathbf{b}_f$ for the forget gate in LSTM , which is initialized to 1.0 (Jozefowicz et al., 2015).

**Optimization Algorithm** Parameter optimization is performed with the Adam optimizer (Kingma and Ba, 2014) with $\beta 1 = \beta 2 = 0.9$. We choose an initial learning rate of $\eta_0 = 0.002$. The learning rate $\eta$ was adapted using a schedule $S = [e_1, e_2, \ldots, e_s]$, in which the learning rate $\eta$ is annealed by multiplying a fixed decay rate $\rho = 0.5$ after $e_i \in S$ epochs respectively. We used $S = [10, 30, 50, 70, 100]$ and trained all networks for a total of 120 epochs. While the Adam

47

|        | English |  |  |  | Chinese |  |  |  | German |  |  |  |
|--------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|
|        | Dev |  | Test |  | Dev |  | Test |  | Dev |  | Test |  |
| Model  | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS |
| Basic  | 94.51 | 92.23 | 94.62 | 92.54 | 84.33 | 81.65 | 84.35 | 81.63 | 90.46 | 87.77 | 90.69 | 88.42 |
| +Char  | 94.74 | 92.55 | 94.73 | 92.75 | 85.07 | 82.63 | 85.24 | 82.46 | 92.16 | 89.82 | 92.24 | 90.18 |
| +POS   | 94.71 | 92.60 | 94.83 | 92.96 | 88.98 | 87.55 | 89.05 | 87.74 | 91.94 | 89.51 | 92.19 | 90.05 |
| Full   | 94.77 | 92.66 | 94.88 | 92.98 | 88.51 | 87.16 | 88.79 | 87.47 | 92.37 | 90.09 | 92.58 | 90.54 |

Table 4.1: Parsing performance (UAS and LAS) of different versions of our model on both the development and test sets for three languages.

optimizer automatically adjusts the global learning rate according to past gradient magnitudes, we find that this additional decay consistently improves model performance across all settings and languages. To reduce the effects of "gradient exploding", we use a gradient clipping of 5.0 (Pascanu et al., 2013). We explored other optimization algorithms such as stochastic gradient descent (SGD) with momentum, AdaDelta (Zeiler, 2012), or RMSProp (Dauphin et al., 2015), but none of them meaningfully improve upon Adam with learning rate annealing in our preliminary experiments. Meanwhile, Adam significantly accelerates the training procedure.

**Dropout Training.** To mitigate overfitting, we apply the dropout method (Srivastava et al., 2014b; Ma et al., 2017) to regularize our model. As shown in 4.1, we apply dropout on character embeddings before inputting to CNN, and on the input, hidden and output vectors of BLSTM. We apply dropout rate of 0.15 to all the embeddings. For BLSTM, we use the recurrent dropout (Gal and Ghahramani, 2016) with 0.25 dropout rate between hidden states and 0.33 between layers. We found that the model using the new recurrent dropout converged much faster than standard dropout, while achiving similar performance.

### 4.1.4 Experiments Setup

We evaluate our neural probabilistic parser on the same data setup as Kuncoro et al. (2016), namely the English Penn Treebank (PTB version 3.0) (Marcus et al., 1993), the Penn Chinese Treebank (CTB version 5.1) (Xue et al., 2002), and the German CoNLL 2009 corpus (Hajič et al.,

2009). Following previous work, all experiments are evaluated on the metrics of unlabeled and labeled attachment score (UAS and LAS).

## 4.1.5   Main Results

We first construct experiments to dissect the effectiveness of each input information (embeddings) of our neural network architecture by ablation studies. We compare the performance of four versions of our model with different inputs — Basic, +POS, +Char and Full — where the Basic model utilizes only the pretrained word embeddings as inputs, while the +POS and +Char models augments the basic one with POS embedding and character information, respectively. The Full model includes all the three information as inputs. According to the results shown in Table 4.1, +Char model obtains better performance than the Basic model on all the three languages, showing that character-level representations are important for dependency parsing. Second, on English and German, +Char and +POS achieves comparable performance, while on Chinese +POS significantly outperforms +Char model. Finally, it should be noted that the Full model achieves the best accuracy among the four models on English and German, but on Chinese +POS obtains the best. We guess that the POS information is more useful for Chinese than English and German.

|  | **Dev** | | **Test** | |
|---|---|---|---|---|
|  | UAS | LAS | UAS | LAS |
| cross-entropy | 94.10 | 91.52 | 93.77 | 91.57 |
| global-likelihood | 94.77 | 92.66 | 94.88 | 92.98 |

Table 4.2: Parsing performance on PTB with different training objective functions.

Table 4.2 gives the performance on PTB of the parsers trained with two different objective functions — the cross-entropy objective of each word, and our objective based on likelihood for an entire tree. The parser with global likelihood objective outperforms the one with simple cross-entropy objective, demonstrating the effectiveness of the global structured objective.

|  | **English** | | **Chinese** | | **German** | |
| --- | --- | --- | --- | --- | --- | --- |
| **System** | UAS | LAS | UAS | LAS | UAS | LAS |
| Bohnet and Nivre (2012) | – | – | 87.3 | 85.9 | 91.4 | 89.4 |
| Chen and Manning (2014) | 91.8 | 89.6 | 83.9 | 82.4 | – | – |
| Ballesteros et al. (2015) | 91.6 | 89.4 | 85.3 | 83.7 | 88.8 | 86.1 |
| Dyer et al. (2015) | 93.1 | 90.9 | 87.2 | 85.7 | – | – |
| Kiperwasser and Goldberg (2016): graph | 93.1 | 91.0 | 86.6 | 85.1 | – | – |
| Ballesteros et al. (2016) | 93.6 | 91.4 | 87.7 | 86.2 | – | – |
| Wang and Chang (2016) | 94.1 | 91.8 | 87.6 | 86.2 | – | – |
| Zhang et al. (2016) | 94.1 | 91.9 | 87.8 | 86.2 | – | – |
| Cheng et al. (2016) | 94.1 | 91.5 | 88.1 | 85.7 | – | – |
| Andor et al. (2016) | 94.6 | 92.8 | – | – | 90.9 | 89.2 |
| Kuncoro et al. (2016) | 94.3 | 92.1 | 88.9 | 87.3 | 91.6 | 89.2 |
| Dozat and Manning (2017) | **95.7** | **94.1** | **89.3** | **88.2** | **93.5** | **91.4** |
| This work: Basic | 94.6 | 92.5 | 84.4 | 81.6 | 90.7 | 88.4 |
| This work: +Char | 94.7 | 92.8 | 85.2 | 82.5 | 92.2 | 90.2 |
| This work: +POS | 94.8 | 93.0 | 89.1 | 87.7 | 92.2 | 90.1 |
| This work: Full | 94.9 | 93.0 | 88.8 | 87.5 | 92.6 | 90.5 |

Table 4.3: UAS and LAS of four versions of our model on test sets for three languages, together with top-performance parsing systems.

## 4.1.6 Comparison with Previous Work

Table 4.3 illustrates the results of the four versions of our model on the three languages, together with twelve previous top-performance systems for comparison. Our Full model significantly outperforms the graph-based parser proposed in Kiperwasser and Goldberg (2016) which used similar neural network architecture for representation learning. Moreover, our model achieves better results than the parser distillation method (Kuncoro et al., 2016) on all the three languages. The results of our parser are slightly worse than the scores reported in Dozat and Manning (2017). One possible reason is that for labeled dependency parsing Dozat and Manning (2017) used different vectors for different dependency labels to represent each word.

| | Turbo | Tensor | RGB | In-Out | Bi-Att | +POS | Full | Best |
|---|---|---|---|---|---|---|---|---|
| | UAS | UAS | UAS | UAS [LAS] | UAS [LAS] | UAS [LAS] | UAS [LAS] | UAS |
| ar | 79.64 | 79.95 | 80.24 | 79.60 [67.09] | 80.34 [68.58] | 80.05 [67.80] | 80.80 [**69.40**] | **81.12** |
| bg | 93.10 | 93.50 | 93.72 | 92.68 [87.79] | 93.96 [89.55] | 93.66 [89.79] | **94.28** [**90.60**] | 94.02 |
| zh | 89.98 | 92.68 | 93.04 | 92.58 [88.51] | – | **93.44** [90.04] | 93.40 [**90.10**] | 93.04 |
| cs | 90.32 | 90.50 | 90.77 | 88.01 [79.31] | 91.16 [85.14] | 91.04 [85.82] | **91.18** [**85.92**] | 91.16 |
| da | 91.48 | 91.39 | 91.86 | 91.44 [85.55] | 91.56 [85.53] | 91.52 [86.57] | 91.86 [**87.07**] | **92.00** |
| nl | 86.19 | 86.41 | 87.39 | 84.45 [80.31] | 87.15 [82.41] | 87.41 [84.17] | **87.85** [**84.82**] | 87.39 |
| en | 93.22 | 93.02 | 93.25 | 92.45 [89.43] | – | 94.43 [92.31] | **94.66** [**92.52**] | 93.25 |
| de | 92.41 | 91.97 | 92.67 | 90.79 [87.74] | 92.71 [89.80] | 93.53 [91.55] | **93.62** [**91.90**] | 92.71 |
| ja | 93.52 | 93.71 | 93.56 | 93.54 [91.80] | 93.44 [90.67] | 93.82 [92.34] | **94.02** [**92.60**] | 93.80 |
| pt | 92.69 | 91.92 | 92.36 | 91.54 [87.68] | 92.77 [88.44] | 92.59 [**89.12**] | 92.71 [88.92] | **93.03** |
| sl | 86.01 | 86.24 | 86.72 | 84.39 [73.74] | 86.01 [75.90] | 85.73 [76.48] | 86.73 [**77.56**] | **87.06** |
| es | 85.59 | 88.00 | 88.75 | 86.44 [83.29] | 88.74 [84.03] | 88.58 [85.03] | **89.20** [**85.77**] | 88.75 |
| sv | 91.14 | 91.00 | 91.08 | 89.94 [83.09] | 90.50 [84.05] | 90.89 [86.58] | 91.22 [**86.92**] | **91.85** |
| tr | 76.90 | 76.84 | 76.68 | 75.32 [60.39] | **78.43** [**66.16**] | 75.88 [61.72] | 77.71 [65.81] | 78.43 |
| av | 88.73 | 89.08 | 89.44 | 88.08 [81.84] | – | 89.47 [84.24] | 89.95 [84.99] | 89.83 |

Table 4.4: UAS and LAS on 14 treebanks from CoNLL shared tasks, together with several state-of-the-art parsers. "Best Published" includes the most accurate parsers in term of UAS among Koo et al. (2010), Martins et al. (2011), Martins et al. (2013), Lei et al. (2014), Zhang et al. (2014), Zhang and McDonald (2014), Pitler and McDonald (2015), and Ma and Hovy (2015).

## 4.1.7 Experiments on CoNLL Treebanks

**Datasets.** To make a thorough empirical comparison with previous studies, we also evaluate our system on treebanks from CoNLL shared task on dependency parsing — the English treebank from CoNLL-2008 shared task (Surdeanu et al., 2008) and all 13 treebanks from CoNLL-2006 shared task (Buchholz and Marsi, 2006). For the treebanks from CoNLL-2006 shared task, following Cheng et al. (2016), we randomly select 5% of the training data as the development set. UAS and LAS are evaluated using the official scorer[1] of CoNLL-2006 shared task.

[1]http://ilk.uvt.nl/conll/software.html

**Baselines.** We compare our model with the third-order Turbo parser (Martins et al., 2013), the low-rank tensor based model (Tensor) (Lei et al., 2014), the randomized greedy inference based (RGB) model (Zhang et al., 2014), the labeled dependency parser with inner-to-outer greedy decoding algorithm (In-Out) (Ma and Hovy, 2015), and the bi-direction attention based parser (Bi-Att) (Cheng et al., 2016). We also compare our parser against the best published results for individual languages. This comparison includes four additional systems: Koo et al. (2010), Martins et al. (2011), Zhang and McDonald (2014) and Pitler and McDonald (2015).

**Results.** Table 4.4 summarizes the results of our model, along with the state-of-the-art baselines. On average across 14 languages, our approach significantly outperforms all the baseline systems. It should be noted that the average UAS of our parser over the 14 languages is better than that of the "best published", which are from different systems that achieved best results for different languages.

For individual languages, our parser achieves state-of-the-art performance on both UAS and LAS on 8 languages — Bulgarian, Chinese, Czech, Dutch, English, German, Japanese and Spanish. On Arabic, Danish, Portuguese, Slovene and Swedish, our parser obtains the best LAS. Another interesting observation is that the Full model outperforms the +POS model on 13 languages. The only exception is Chinese, which matches the observation in Section 4.1.5.

## 4.2 Discussions

In this chapter, we proposed a neural probabilistic model for non-projective dependency parsing, which combined the BLSTM-CNNs architecture for representation learning with a probabilistic structured output layer on top. Experimental results on 17 treebanks across 14 languages show that our parser significantly improves the accuracy of both dependency structures (UAS) and edge labels (LAS), over several previously state-of-the-art systems.

In the next chapter, we will consider neural networks for transition-based dependency parsing, where we introduce our stack-pointer network.

# Chapter 5

# Stack-Pointer Networks for Dependency Parsing

In the last chapter, we showed that incorporating this global search algorithm with distributed representations learned from neural networks, neural graph-based parsers (Kiperwasser and Goldberg, 2016; Wang and Chang, 2016; Kuncoro et al., 2016; Dozat and Manning, 2017) have achieved the state-of-the-art accuracies on a number of treebanks in different languages. Nevertheless, these models, while accurate, are usually slow (e.g. decoding is $O(n^3)$ time complexity for first-order models McDonald et al. (2005a,b) and higher polynomials for higher-order models (McDonald and Pereira, 2006; Koo and Collins, 2010; Ma and Zhao, 2012b,a)).

Transition-based dependency parsers, on the other hand, read words sequentially (commonly from left-to-right) and build dependency trees incrementally by making series of multiple choice decisions. A classifier is trained to score the possible decisions at each state of the process and guide the parsing process. The advantage of this formalism is that the number of operations required to build any projective parse tree is linear with respect to the length of the sentence. The challenge, however, is that the decision made at each step is based on local information, leading to error propagation and worse performance compared to graph-based parsers on root and long dependencies (McDonald and Nivre, 2011). Previous studies have explored solutions to address this challenge. Stack LSTMs (Dyer et al., 2015; Ballesteros et al., 2015, 2016) are capable of learning representations of the parser state that are sensitive to the complete contents of

the parser's state. Andor et al. (2016) proposed a globally normalized transition model to replace the locally normalized classifier. However, the parsing accuracy is still behind state-of-the-art graph-based parsers (Dozat and Manning, 2017).

In this chapter, we propose a novel neural network architecture for dependency parsing, *stack-pointer networks* (**STACKPTR**). STACKPTR is a transition-based architecture, with the corresponding asymptotic efficiency, but still maintains a global view of the sentence that proves essential for achieving competitive accuracy. Our STACKPTR parser has a pointer network (Vinyals et al., 2015) as its backbone, and is equipped with an internal stack to maintain the order of head words in tree structures. The STACKPTR parser performs parsing in an incremental, top-down, depth-first fashion; at each step, it generates an arc by assigning a child for the head word at the top of the internal stack. This architecture makes it possible to capture information from the whole sentence and all the previously derived subtrees, while maintaining a number of parsing steps linear in the sentence length.

We evaluate our parser on 29 treebanks across 20 languages and different dependency annotation schemas, and achieve state-of-the-art performance on 21 of them. The contributions of this work are summarized as follows:

(i) We propose a neural network architecture for dependency parsing that is simple, effective, and efficient.

(ii) Empirical evaluations on benchmark datasets over 20 languages show that our method achieves state-of-the-art performance on 21 different treebanks.

(iii) Comprehensive error analysis is conducted to compare the proposed method to a strong graph-based baseline using biaffine attention (Dozat and Manning, 2017).

Source code for the implementation is publicly available at `https://github.com/XuezheMax/NeuroNLP2`.

## 5.1 Background

### 5.1.1 Notations

Dependency trees represent syntactic relationships between words in the sentences through labeled directed edges between head words and their dependents. Figure 5.1 (a) shows a dependency tree for the sentence, "But there were no buyers". In this paper, we will use the following notation:

**Input**: $\mathbf{x} = \{w_1, \ldots, w_n\}$ represents a generic sentence, where $w_i$ is the $i$th word.

**Output**: $\mathbf{y} = \{p_1, p_2, \cdots, p_k\}$ represents a generic (possibly non-projective) dependency tree, where each path $p_i = \$, w_{i,1}, w_{i,2}, \cdots, w_{i,l_i}$ is a sequence of words from the root to a leaf. "$\$$" is an universal virtual root that is added to each tree.

**Stack**: $\sigma$ denotes a stack configuration, which is a sequence of words. We use $\sigma|w$ to represent a stack configuration that pushes word $w$ into the stack $\sigma$.

**Children**: $\mathrm{ch}(w_i)$ denotes the list of all the children (modifiers) of word $w_i$.

### 5.1.2 Pointer Networks

Pointer Networks (PTR-NET) (Vinyals et al., 2015) are a variety of neural network capable of learning the conditional probability of an output sequence with elements that are discrete tokens corresponding to positions in an input sequence. This model cannot be trivially expressed by standard sequence-to-sequence networks (Sutskever et al., 2014) due to the variable number of input positions in each sentence. PTR-NET solves the problem by using attention (Bahdanau et al., 2015; Luong et al., 2015) as a pointer to select a member of the input sequence as the output.

Formally, the words of the sentence $\mathbf{x}$ are fed one-by-one into the encoder (a multiple-layer bi-directional RNN), producing a sequence of *encoder hidden states* $s_i$. At each time step $t$, the decoder (a uni-directional RNN) receives the input from last step and outputs *decoder hidden state* $h_t$. The *attention vector* $a^t$ is calculated as follows:

$$
\begin{aligned}
e_i^t &= score(h_t, s_i) \\
a^t &= softmax(e^t)
\end{aligned}
\tag{5.1}
$$

where $score(\cdot, \cdot)$ is the *attention scoring function*, which has several variations such as dot-product,

(a)                                        (b)

Figure 5.1: Neural architecture for the STACKPTR network, together with the decoding procedure of an example sentence. The BiRNN of the encoder is elided for brevity. For the inputs of decoder at each time step, vectors in red and blue boxes indicate the sibling and grandparent.

concatenation, and biaffine (Luong et al., 2015). PTR-NET regards the attention vector $a^t$ as a probability distribution over the source words, i.e. it uses $a^t_i$ as pointers to select input elements.

## 5.2 Stack-Pointer Networks

### 5.2.1 Overview

Similarly to PTR-NET, STACKPTR first reads the whole sentence and encodes each word into the encoder hidden state $s_i$. The internal stack $\sigma$ is always initialized with the root symbol \$. At each time step $t$, the decoder receives the input vector corresponding to the top element of the stack $\sigma$ (the head word $w_p$ where $p$ is the word index), generates the hidden state $h_t$, and computes the attention vector $a^t$ using Eq. (5.1). The parser chooses a specific position $c$ according to the attention scores in $a^t$ to generate a new dependency arc $(w_h, w_c)$ by selecting $w_c$ as a child of $w_h$. Then the parser pushes $w_c$ onto the stack, i.e. $\sigma \rightarrow \sigma|w_c$, and goes to the next step. At one step if the parser points $w_h$ to itself, i.e. $c = h$, it indicates that all children of the head word $w_h$ have

already been selected. Then the parser goes to the next step by popping $w_h$ out of $\sigma$.

At test time, in order to guarantee a valid dependency tree containing all the words in the input sentences exactly once, the decoder maintains a list of "available" words. At each decoding step, the parser selects a child for the current head word, and removes the child from the list of available words to make sure that it cannot be selected as a child of other head words.

For head words with multiple children, it is possible that there is more than one valid selection for each time step. In order to define a deterministic decoding process to make sure that there is only one ground-truth choice at each step (which is necessary for simple maximum likelihood estimation), a predefined order for each $\mathrm{ch}(w_i)$ needs to be introduced. The predefined order of children can have different alternatives, such as left-to-right or inside-out[1]. In this paper, we adopt the inside-out order[2] since it enables us to utilize second-order *sibling* information, which has been proven beneficial for parsing performance (McDonald and Pereira, 2006; Koo and Collins, 2010) (see § 5.2.4 for details). Figure 5.1 (b) depicts the architecture of STACKPTR and the decoding procedure for the example sentence in Figure 5.1 (a).

## 5.2.2 Encoder

The encoder of our parsing model is based on the bi-directional LSTM-CNN architecture (BLSTM-CNNs) (Ma and Hovy, 2016) where CNNs encode character-level information of a word into its character-level representation and BLSTM models context information of each word. Formally, for each word, the CNN, with character embeddings as inputs, encodes the character-level representation. Then the character-level representation vector is concatenated with the word embedding vector to feed into the BLSTM network. To enrich word-level information, we also use POS embeddings. Finally, the encoder outputs a sequence of hidden states $s_i$.

## 5.2.3 Decoder

The decoder for our parser is a uni-directional LSTM. Different from previous work (Bahdanau et al., 2015; Vinyals et al., 2015) which uses word embeddings of the previous word as the input

---

[1] Order the children by the distances to the head word on the left side, then the right side.
[2] We also tried left-to-right order which obtained worse parsing accuracy than inside-out.

to the decoder, our decoder receives the encoder hidden state vector ($s_i$) of the top element in the stack $\sigma$ (see Figure 5.1 (b)). Compared to word embeddings, the encoder hidden states contain more contextual information, benefiting both the training and decoding procedures. The decoder produces a sequence of decoder hidden states $h_i$, one for each decoding step.

### 5.2.4 Higher-order Information

As mentioned before, our parser is capable of utilizing higher-order information. In this paper, we incorporate two kinds of higher-order structures — *grandparent* and *sibling*. A sibling structure is a head word with two successive modifiers, and a grandparent structure is a pair of dependencies connected head-to-tail:



To utilize higher-order information, the decoder's input at each step is the sum of the encoder hidden states of three words:

$$\beta_t = s_h + s_g + s_s$$

where $\beta_t$ is the input vector of decoder at time $t$ and $h, g, s$ are the indices of the head word and its grandparent and sibling, respectively. Figure 5.1 (b) illustrates the details. Here we use the element-wise sum operation instead of concatenation because it does not increase the dimension of the input vector $\beta_t$, thus introducing no additional model parameters.

### 5.2.5 Biaffine Attention Mechanism

For attention score function (Eq. (5.1)), we adopt the biaffine attention mechanism (Luong et al., 2015; Dozat and Manning, 2017):

$$e_i^t = h_t^T \mathbf{W} s_i + \mathbf{U}^T h_t + \mathbf{V}^T s_i + \mathbf{b}$$

where $\mathbf{W}, \mathbf{U}, \mathbf{V}, \mathbf{b}$ are parameters, denoting the weight matrix of the bi-linear term, the two weight vectors of the linear terms, and the bias vector.

As discussed in Dozat and Manning (2017), applying a multilayer perceptron (MLP) to the output vectors of the BLSTM before the score function can both reduce the dimensionality and overfitting of the model. We follow this work by using a one-layer perceptron to $s_i$ and $h_i$ with ELU Clevert et al. (2015) as its activation function. Similarly, the dependency label classifier also uses a biaffine function to score each label, given the head word vector $h_t$ and child vector $s_i$ as inputs. Again, we use MLPs to transform $h_t$ and $s_i$ before feeding them into the classifier.

## 5.2.6 Training Objectives

The STACKPTR parser is trained to optimize the probability of the dependency trees given sentences: $P_\theta(\mathbf{y}|\mathbf{x})$, which can be factorized as:

$$P_\theta(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{k} P_\theta(p_i|p_{<i},\mathbf{x}) = \prod_{i=1}^{k}\prod_{j=1}^{l_i} P_\theta(c_{i,j}|c_{i,<j}, p_{<i}, \mathbf{x}), \tag{5.2}$$

where $\theta$ represents model parameters. $p_{<i}$ denotes the preceding paths that have already been generated. $c_{i,j}$ represents the $j$th word in $p_i$ and $c_{i,<j}$ denotes all the proceeding words on the path $p_i$. Thus, the STACKPTR parser is an autoregressive model, like sequence-to-sequence models, but it factors the distribution according to a top-down tree structure as opposed to a left-to-right chain. We define $P_\theta(c_{i,j}|c_{i,<j}, p_{<i}, \mathbf{x}) = a^t$, where attention vector $a^t$ (of dimension $n$) is used as the distribution over the indices of words in a sentence.

**Arc Prediction**    Our parser is trained by optimizing the conditional likelihood in Eq (5.2), which is implemented as the cross-entropy loss.

**Label Prediction**    We train a separated multi-class classifier in parallel to predict the dependency labels. Following Dozat and Manning (2017), the classifier takes the information of the head word and its child as features. The label classifier is trained simultaneously with the parser by optimizing the sum of their objectives.

## 5.2.7 Discussion

**Time Complexity.**    The number of decoding steps to build a parse tree for a sentence of length $n$ is $2n - 1$, linear in $n$. Together with the attention mechanism (at each step, we need to compute

the attention vector $a^t$, whose runtime is $O(n)$), the time complexity of decoding algorithm is $O(n^2)$, which is more efficient than graph-based parsers that have $O(n^3)$ or worse complexity when using dynamic programming or maximum spanning tree (MST) decoding algorithms.

**Top-down Parsing.** When humans comprehend a natural language sentence, they arguably do it in an incremental, left-to-right manner. However, when humans consciously annotate a sentence with syntactic structure, they rarely ever process in fixed left-to-right order. Rather, they start by reading the whole sentence, then seeking the main predicates, jumping back-and-forth over the sentence and recursively proceeding to the sub-tree structures governed by certain head words. Our parser follows a similar kind of annotation process: starting from reading the whole sentence, and processing in a top-down manner by finding the main predicates first and only then search for sub-trees governed by them. When making latter decisions, the parser has access to the entire structure built in earlier steps.

### 5.2.8 Implementation Details

**Pre-trained Word Embeddings.** For all the parsing models in different languages, we initialize word vectors with pretrained word embeddings. For Chinese, Dutch, English, German and Spanish, we use the structured-skipgram Ling et al. (2015) embeddings. For other languages we use Polyglot embeddings Al-Rfou et al. (2013).

**Optimization.** Parameter optimization is performed with the Adam optimizer Kingma and Ba (2014) with $\beta_1 = \beta_2 = 0.9$. We choose an initial learning rate of $\eta_0 = 0.001$. The learning rate $\eta$ is annealed by multiplying a fixed decay rate $\rho = 0.75$ when parsing performance stops increasing on validation sets. To reduce the effects of "gradient exploding", we use gradient clipping of $5.0$ Pascanu et al. (2013).

**Dropout Training.** To mitigate overfitting, we apply dropout Srivastava et al. (2014b); Ma et al. (2017). For BLSTM, we use recurrent dropout Gal and Ghahramani (2016) with a drop rate of 0.33 between hidden states and 0.33 between layers. Following Dozat and Manning (2017), we also use embedding dropout with a rate of 0.33 on all word, character, and POS embeddings.

## 5.3 Experiments

### 5.3.1 Setup

We evaluate our STACKPTR parser mainly on three treebanks: the English Penn Treebank (PTB version 3.0) (Marcus et al., 1993), the Penn Chinese Treebank (CTB version 5.1) Xue et al. (2002), and the German CoNLL 2009 corpus Hajič et al. (2009). We use the same experimental settings as Kuncoro et al. (2016).

To make a thorough empirical comparison with previous studies, we also evaluate our system on treebanks from CoNLL shared task and the Universal Dependency (UD) Treebanks[3]. For the CoNLL Treebanks, we use the English treebank from CoNLL-2008 shared task Surdeanu et al. (2008) and all 13 treebanks from CoNLL-2006 shared task Buchholz and Marsi (2006). The experimental settings are the same as Ma and Hovy (2015). For UD Treebanks, we select 12 languages. The details of the treebanks and experimental settings are in § 5.3.6 and Appendix A.1.2.

**Evaluation Metrics**    Parsing performance is measured with five metrics: unlabeled attachment score (UAS), labeled attachment score (LAS), unlabeled complete match (UCM), labeled complete match (LCM), and root accuracy (RA). Following previous work (Kuncoro et al., 2016; Dozat and Manning, 2017), we report results excluding punctuations for Chinese and English. For each experiment, we report the mean values with corresponding standard deviations over 5 repetitions.

**Baseline**    For fair comparison of the parsing performance, we re-implemented the graph-based Deep Biaffine (BIAF) parser (Dozat and Manning, 2017), which achieved state-of-the-art results on a wide range of languages. Our re-implementation adds character-level information using the same LSTM-CNN encoder as our model (§ 5.2.2) to the original BIAF model, which boosts its performance on all languages.

### 5.3.2 Main Results

We first conduct experiments to demonstrate the effectiveness of our neural architecture by comparing with the strong baseline BIAF. We compare the performance of four variations of

---

[3]http://universaldependencies.org/

Figure 5.2: Parsing performance of different variations of our model on the test sets for three languages, together with baseline BIAF. For each of our STACKPTR models, we perform decoding with beam size equal to 1 and 10. The improvements of decoding with beam size 10 over 1 are presented by stacked bars with light colors.

our model with different decoder inputs — Org, +gpar, +sib and Full — where the Org model utilizes only the encoder hidden states of head words, while the +gpar and +sib models augments the original one with grandparent and sibling information, respectively. The Full model includes all the three information as inputs.

Figure 5.2 illustrates the performance (five metrics) of different variations of our STACKPTR parser together with the results of baseline BIAF re-implemented by us, on the test sets of the three languages. On UAS and LAS, the Full variation of STACKPTR with decoding beam size 10 outperforms BIAF on Chinese, and obtains competitive performance on English and German. An interesting observation is that the Full model achieves the best accuracy on English and Chinese, while performs slightly worse than +sib on German. This shows that the importance of higher-order information varies in languages. On LCM and UCM, STACKPTR significantly

62

| System | | English | | Chinese | | German | |
|---|---|---|---|---|---|---|---|
| | | UAS | LAS | UAS | LAS | UAS | LAS |
| Chen and Manning (2014) | T | 91.8 | 89.6 | 83.9 | 82.4 | – | – |
| Ballesteros et al. (2015) | T | 91.63 | 89.44 | 85.30 | 83.72 | 88.83 | 86.10 |
| Dyer et al. (2015) | T | 93.1 | 90.9 | 87.2 | 85.7 | – | – |
| Bohnet and Nivre (2012) | T | 93.33 | 91.22 | 87.3 | 85.9 | 91.4 | 89.4 |
| Ballesteros et al. (2016) | T | 93.56 | 91.42 | 87.65 | 86.21 | – | – |
| Kiperwasser and Goldberg (2016) | T | 93.9 | 91.9 | 87.6 | 86.1 | – | – |
| Weiss et al. (2015b) | T | 94.26 | 92.41 | – | – | – | – |
| Andor et al. (2016) | T | 94.61 | 92.79 | – | – | 90.91 | 89.15 |
| Kiperwasser and Goldberg (2016) | G | 93.1 | 91.0 | 86.6 | 85.1 | – | – |
| Wang and Chang (2016) | G | 94.08 | 91.82 | 87.55 | 86.23 | – | – |
| Cheng et al. (2016) | G | 94.10 | 91.49 | 88.1 | 85.7 | – | – |
| Kuncoro et al. (2016) | G | 94.26 | 92.06 | 88.87 | 87.30 | 91.60 | 89.24 |
| Ma and Hovy (2017) | G | 94.88 | 92.98 | 89.05 | 87.74 | 92.58 | 90.54 |
| BIAF: Dozat and Manning (2017) | G | 95.74 | 94.08 | 89.30 | 88.23 | 93.46 | 91.44 |
| BIAF: re-impl | G | 95.84 | **94.21** | 90.43 | 89.14 | **93.85** | **92.32** |
| STACKPTR: Org | T | 95.77 | 94.12 | 90.48 | 89.19 | 93.59 | 92.06 |
| STACKPTR: +gpar | T | 95.78 | 94.12 | 90.49 | 89.19 | 93.65 | 92.12 |
| STACKPTR: +sib | T | 95.85 | 94.18 | 90.43 | 89.15 | 93.76 | 92.21 |
| STACKPTR: Full | T | **95.87** | 94.19 | **90.59** | **89.29** | 93.65 | 92.11 |

Table 5.1: UAS and LAS of four versions of our model on test sets for three languages, together with top-performing parsing systems. "T" and "G" indicate transition- and graph-based models, respectively. For BIAF, we provide the original results reported in Dozat and Manning (2017) and our re-implementation.

outperforms BIAF on all languages, showing the superiority of our parser on complete sentence parsing. The results of our parser on RA are slightly worse than BIAF. More details of results are provided in Appendix A.1.2.

### 5.3.3 Comparison with Previous Work

Table 5.1 illustrates the UAS and LAS of the four versions of our model (with decoding beam size 10) on the three treebanks, together with previous top-performing systems for comparison. Note that the results of STACKPTR and our re-implementation of BIAF are the average of 5 repetitions instead of a single run. Our Full model significantly outperforms all the transition-

Figure 5.3: Performance of BIAF and STACKPTR parsers relative to length and graph factors.

based parsers on all three languages, and achieves better results than most graph-based parsers. Our re-implementation of BIAF obtains better performance than the original one in Dozat and Manning (2017), demonstrating the effectiveness of the character-level information. Our model achieves state-of-the-art performance on both UAS and LAS on Chinese, and best UAS on English. On German, the performance is competitive with BIAF, and significantly better than other models.

### 5.3.4 Error Analysis

In this section, we characterize the errors made by BIAF and STACKPTR by presenting a number of experiments that relate parsing errors to a set of linguistic and structural properties. For simplicity, we follow McDonald and Nivre (2011) and report labeled parsing metrics (either accuracy, precision, or recall) for all experiments.

**Length and Graph Factors**

Following McDonald and Nivre (2011), we analyze parsing errors related to structural factors.

**Sentence Length.** Figure 5.3 (a) shows the accuracy of both parsing models relative to sentence lengths. Consistent with the analysis in McDonald and Nivre (2011), STACKPTR tends to perform

| POS | UAS | LAS | UCM | LCM |
|------|-----------|-----------|-----------|-----------|
| Gold | 96.12±0.03 | 95.06±0.05 | 62.22±0.33 | 55.74±0.44 |
| Pred | 95.87±0.04 | 94.19±0.04 | 61.43±0.49 | 49.68±0.47 |
| None | 95.90±0.05 | 94.21±0.04 | 61.58±0.39 | 49.87±0.46 |

Table 5.2: Parsing performance on the test data of PTB with different versions of POS tags.

better on shorter sentences, which make fewer parsing decisions, significantly reducing the chance of error propagation.

**Dependency Length.**    Figure 5.3 (b) measures the precision and recall relative to dependency lengths. While the graph-based BIAF parser still performs better for longer dependency arcs and transition-based STACKPTR parser does better for shorter ones, the gap between the two systems is marginal, much smaller than that shown in McDonald and Nivre (2011). One possible reason is that, unlike traditional transition-based parsers that scan the sentence from left to right, STACKPTR processes in a top-down manner, unnecessarily creating shorter dependencies first.

**Root Distance.**    Figure 5.3 (c) plots the precision and recall of each system for arcs of varying distance to the root. Different from the observation in McDonald and Nivre (2011), STACKPTR does not show an obvious advantage on the precision for arcs further away from the root. Furthermore, the STACKPTR parser does not have the tendency to over-predict root modifiers reported in McDonald and Nivre (2011). This behavior can be explained using the same reasoning as above: the fact that arcs further away from the root are usually constructed early in the parsing algorithm of traditional transition-based parsers is not true for the STACKPTR parser.

**Effect of POS Embedding**

The only prerequisite information that our parsing model relies on is POS tags. With the goal of achieving an end-to-end parser, we explore the effect of POS tags on parsing performance. We run experiments on PTB using our STACKPTR parser with gold-standard and predicted POS tags, and without tags, respectively. STACKPTR in these experiments is the Full model with beam=10.

Table 5.2 gives results of the parsers with different versions of POS tags on the test data of PTB. The parser with gold-standard POS tags significantly outperforms the other two parsers,

|  | Bi-Att | NeuroMST | BIAF | STACKPTR | Best Published | |
|---|---|---|---|---|---|---|
|  | UAS [LAS] | UAS [LAS] | UAS [LAS] | UAS [LAS] | UAS | LAS |
| ar | 80.34 [68.58] | 80.80 [69.40] | 82.15±0.34 [71.32±0.36] | **83.04±0.29 [72.94±0.31]** | 81.12 | – |
| bg | 93.96 [89.55] | 94.28 [90.60] | 94.62±0.14 [**91.56±0.24**] | **94.66±0.10** [91.40±0.08] | 94.02 | – |
| zh | – | 93.40 [90.10] | **94.05±0.27 [90.89±0.22]** | 93.88±0.24 [90.81±0.55] | 93.04 | – |
| cs | 91.16 [85.14] | 91.18 [85.92] | 92.24±0.22 [87.85±0.21] | **92.83±0.13 [88.75±0.16]** | 91.16 | 85.14 |
| da | 91.56 [85.53] | 91.86 [87.07] | **92.80±0.26 [88.36±0.18]** | 92.08±0.15 [87.29±0.21] | 92.00 | – |
| nl | 87.15 [82.41] | 87.85 [84.82] | 90.07±0.18 [**87.24±0.17**] | **90.10±0.27** [87.05±0.26] | 87.39 | – |
| en | – | 94.66 [92.52] | 95.19±0.05 [93.14±0.05] | **93.25±0.05 [93.17±0.05]** | 93.25 | – |
| de | 92.71 [89.80] | 93.62 [91.90] | 94.52±0.11 [93.06±0.11] | **94.77±0.05 [93.21±0.10]** | 92.71 | 89.80 |
| ja | 93.44 [90.67] | **94.02 [92.60]** | 93.95±0.06 [92.46±0.07] | 93.38±0.08 [91.92±0.16] | 93.80 | – |
| pt | 92.77 [88.44] | 92.71 [88.92] | 93.41±0.08 [89.96±0.24] | **93.57±0.12 [90.07±0.20]** | 93.03 | – |
| sl | 86.01 [75.90] | 86.73 [77.56] | 87.55±0.17 [78.52±0.35] | **87.59±0.36 [78.85±0.53]** | 87.06 | – |
| es | 88.74 [84.03] | 89.20 [85.77] | 90.43±0.13 [87.08±0.14] | **90.87±0.26 [87.80±0.31]** | 88.75 | 84.03 |
| sv | 90.50 [84.05] | 91.22 [86.92] | 92.22±0.15 [88.44±0.17] | **92.49±0.21 [89.01±0.22]** | 91.85 | 85.26 |
| tr | 78.43 [66.16] | 77.71 [65.81] | **79.84±0.23 [68.63±0.29]** | 79.56±0.22 [68.03±0.15] | 78.43 | 66.16 |

Table 5.3: UAS and LAS on 14 treebanks from CoNLL shared tasks, together with several state-of-the-art parsers. Bi-Att is the bi-directional attention based parser (Cheng et al., 2016), and NeuroMST is the neural MST parser (Ma and Hovy, 2017). "Best Published" includes the most accurate parsers in term of UAS among Koo et al. (2010), Martins et al. (2011), Martins et al. (2013), Lei et al. (2014), Zhang et al. (2014), Zhang and McDonald (2014), Pitler and McDonald (2015), and Cheng et al. (2016).

showing that dependency parsers can still benefit from accurate POS information. The parser with predicted (imperfect) POS tags, however, performs even slightly worse than the parser without using POS tags. It illustrates that an end-to-end parser that doesn't rely on POS information can obtain competitive (or even better) performance than parsers using imperfect predicted POS tags, even if the POS tagger is relative high accuracy (accuracy $> 97\%$ in this experiment on PTB).

### 5.3.5 Experiments on CoNLL Treebanks

Table 5.3 summarizes the parsing results of our model on the test sets of 14 treebanks from the CoNLL shared task, along with the state-of-the-art baselines. Along with BIAF, we also list

| | Dev | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| | BIAF | | STACKPTR | | BIAF | | STACKPTR | |
| | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS |
| bg | 93.92±0.13 | 89.05±0.11 | **94.09±0.16** | **89.17±0.14** | 94.30±0.16 | **90.04±0.16** | **94.31±0.06** | 89.96±0.07 |
| ca | 94.21±0.05 | 91.97±0.06 | **94.47±0.02** | **92.51±0.05** | 94.36±0.06 | 92.05±0.07 | **94.47±0.02** | **92.39±0.02** |
| cs | 94.14±0.03 | 90.89±0.04 | **94.33±0.04** | **91.24±0.05** | 94.06±0.04 | 90.60±0.05 | **94.21±0.06** | **90.94±0.07** |
| de | 91.89±0.11 | 88.39±0.17 | **92.26±0.11** | **88.79±0.15** | 90.26±0.19 | 86.11±0.25 | **90.26±0.07** | **86.16±0.01** |
| en | **92.51±0.08** | **90.50±0.07** | 92.47±0.03 | 90.46±0.02 | 91.91±0.17 | 89.82±0.16 | **91.93±0.07** | **89.83±0.06** |
| es | 93.46±0.05 | 91.13±0.07 | **93.54±0.06** | **91.34±0.05** | 93.72±0.07 | 91.33±0.08 | **93.77±0.07** | **91.52±0.07** |
| fr | **95.05±0.04** | **92.76±0.07** | 94.97±0.04 | 92.57±0.06 | 92.62±0.15 | 89.51±0.14 | **92.90±0.20** | **89.88±0.23** |
| it | 94.89±0.12 | 92.58±0.12 | **94.93±0.09** | **92.90±0.10** | **94.75±0.12** | **92.72±0.12** | 94.70±0.07 | 92.55±0.09 |
| nl | 93.39±0.08 | 90.90±0.07 | **93.94±0.11** | **91.67±0.08** | 93.44±0.09 | 91.04±0.06 | **93.98±0.05** | **91.73±0.07** |
| no | 95.44±0.05 | 93.73±0.05 | **95.52±0.08** | **93.80±0.08** | 95.28±0.05 | 93.58±0.05 | **95.33±0.03** | **93.62±0.03** |
| ro | 91.97±0.13 | 85.38±0.03 | **92.06±0.08** | **85.58±0.12** | **91.94±0.07** | **85.61±0.13** | 91.80±0.11 | 85.34±0.21 |
| ru | 93.81±0.05 | 91.85±0.06 | **94.11±0.07** | **92.29±0.10** | 94.40±0.03 | 92.68±0.04 | **94.69±0.04** | **93.07±0.03** |

Table 5.4: UAS and LAS on both the development and test datasets of 12 treebanks from UD Treebanks, together with BIAF for comparison.

the performance of the bi-directional attention based Parser (Bi-Att) (Cheng et al., 2016) and the neural MST parser (NeuroMST) (Ma and Hovy, 2017) for comparison. Our parser achieves state-of-the-art performance on both UAS and LAS on eight languages — Arabic, Czech, English, German, Portuguese, Slovene, Spanish, and Swedish. On Bulgarian and Dutch, our parser obtains the best UAS. On other languages, the performance of our parser is competitive with BIAF, and significantly better than others. The only exception is Japanese, on which NeuroMST obtains the best scores.

### 5.3.6   Experiments on UD Treebanks

For UD Treebanks, we select 12 languages — Bulgarian, Catalan, Czech, Dutch, English, French, German, Italian, Norwegian, Romanian, Russian and Spanish. For all the languages, we adopt the standard training/dev/test splits, and use the universal POS tags (Petrov et al., 2012) provided in each treebank. The statistics of these corpora are provided in Appendix A.2. For evaluation, we

report results excluding punctuation, which is any tokens with POS tags "PUNCT" or "SYM".

Table 5.4 summarizes the results of the STACKPTR parser, along with BIAF for comparison, on both the development and test datasets for each language. First, both BIAF and STACKPTR parsers achieve relatively high parsing accuracies on all the 12 languages — all with UAS are higher than 90%. On nine languages — Catalan, Czech, Dutch, English, French, German, Norwegian, Russian and Spanish — STACKPTR outperforms BIAF for both UAS and LAS. On Bulgarian, STACKPTR achieves slightly better UAS while LAS is slightly worse than BIAF. On Italian and Romanian, BIAF obtains marginally better parsing performance than STACKPTR.

## 5.4 Discussions

In this chapter, we proposed STACKPTR, a transition-based neural network architecture, for dependency parsing. Combining pointer networks with an internal stack to track the status of the top-down, depth-first search in the decoding procedure, the STACKPTR parser is able to capture information from the whole sentence and all the previously derived subtrees, removing the left-to-right restriction in classical transition-based parsers, while maintaining linear parsing steps, w.r.t the length of the sentences. Experimental results on 29 treebanks show the effectiveness of our parser across 20 languages, by achieving state-of-the-art performance on 21 corpora.

# Part II

# Interpretability of Neural Structured Prediction Models

# Chapter 6

# Interpretability of Deep Neural Networks and the Probing Method

The end-to-end training paradigm significantly simplifies the hand-crafted feature engineering process in traditional feature-based machine learning (ML) systems, while giving the neural models flexibility to be optimized towards the ultimate tasks. This simplicity, however, comes at the expense of model interpretability (Shi et al., 2016; Lipton, 2018; Belinkov et al., 2019). Unlike traditional feature-engineered NLP systems whose features, e.g. morphological properties, syntactic categories or semantic relations, are more easily understood by humans, it is more difficult to understand what happens in the internal components of an end-to-end neural network. It is not clear what the role of different components is, how they interact, and what kind of information they learn during the training process. Consequently, such deep neural models are sometimes perceived as "black-box", hindering research efforts and limiting their utility to society (Belinkov, 2018).

The lack of interpretability has major implications for the adoption and further development of AI systems. Gaining a better understanding of these AI systems is necessary not only for improving their design and performance, but also for guaranteeing fairness and accountability them. Thus, the second part of this thesis focuses on model interpretability of deep neural networks. In this chapter, I first review terminologial issues regarding analysis and interpretation in machine learning (Section 6.1). Then I briefly describe the methodological approach of probing,

which is used throughout the following chapters for analyzing deep learning models (Section 6.2), and survey related work on applying probing method to analyze neural networks. Section 6.3 surveys related work on applying probing method to analyze neural networks, together with a brief summary of other analysis methods that have been considered in this literature. In the last section (Section 6.4), we revisit the probing accuracy on reflecting the quality of linguistic properties given representations in a view of prediction. I apply probes to neural dependency paring models to analyze the part-of-speech (POS) information encoded in their internal representations. We argue that, without considering the expressiveness of probing classifiers, the accuracy of probes does not consistently reflect the quality of the encoded information. Based on the experimental results of a case study on probing neural dependency parsers, we propose to interpret performance of probing tasks with two separate metrics, *capacity* and *accessibility*, which are associated with probe expressiveness.

## 6.1   Terminological Issues on Interpretability

As discussed in Belinkov (2018), terms such as interpretability, explainability, transparency, explainable AI (XAI) have been interchangeably used in the context of work on deep learning, and more broadly machine learning and AI. At present there seems to be no consensus on their precise definition and application to the study of AI systems. I think a short review of aspects of terminology is helpful to clarify the questions this thesis considers in the broader work on interpretability in AI. For detailed discussions, please see (Doshi-Velez and Kim, 2017; Lipton, 2018; Belinkov, 2018), as well as the online book by Christoph Molnar[1], for more references.

**Interpretability from explaining decisions.**
Miller (2019) defined interpretability, from the perspective of explanation in social sciences, as "the degree to which an observer can understand the cause of a decision". In this definition, interpretability is the same as *explainability*. While explaining specific model predictions is obviously important in work on deep learning and is recognized as a desideratum for increasing the accountability of machine learning systems (Doshi-Velez et al., 2017), it is different from

---

[1]https://christophm.github.io/interpretable-ml-book/

interpretability on explicitly explaining decisions for given examples. Doshi-Velez and Kim (2017) defined interpretability as "the ability to explain or to present in understandable terms to human", not referring to *decisions*. More relevant work along these lines is mentioned in Belinkov (2018).

**Interpretability from Transparency.**

Lipton (2018) related interpretability *transparency*, which is concerned with how the model works. One important criterion is the level of analysis. Transparency can operate at the level of the entire model (*simulatability*), at the level of individual components (*decomposability*), and at the level of the training algorithm (*algorithmic transparency*).

From a global level, simulatability considers the transparency of the entire model by asking if a person can contemplate it at once. This definition is , however, too conservative and favors simple models with limited expressiveness such as sparse linear models. Similar notion is adopted in (Ribeiro et al., 2016), suggesting that an interpretable model is one that "can be readily presented to the user with visual or textual artifacts".

From a local level, decomposability requires that each part of the model — input, parameter, and calculation – admits an intuitive explanation. This accords with the property of intelligibility as presented in Lou et al. (2012).

From the level of learning algorithms, algorithmic transparency might provide us confidence on what solution the model converge to and/or what the model will behave on unseen data. Most modern deep learning methods lack this sort of algorithmic transparency.

**Interpretability from post-hoc explanations.**

Lipton (2018) also contrasts transparency with *post-hoc explanation*, which is what else can the model tell us. Post hoc interpretability represents a distinct approach to extracting information from learned models. While post hoc interpretations often do not elucidate precisely how a model works, they may nonetheless confer useful information for practitioners and end users of machine learning. Some common approaches to post-hoc explanations include natural language explanations, visualizations of learned representations or models, and probing the internal layers of deep neural networks.

Figure 6.1: Framework of the methodology of probing method.

### 6.1.1 What This Dissertation Is About

From the perspective of methodology, this thesis provides post-hoc explanations for neural linguistic structured prediction models, by adopting probing methods to analyze the internal representations. From the perspective of motivation, this thesis aims to use probing methods to investigate in what way these neural models memorize and process linguiscit information across their hidden representations. It provides a better understanding of the learning properties of different parts and modules in deep learning models (striving for decomposibility and algorithmic transparency, in the sense of (Lipton, 2018)).

## 6.2 Methodology of Probing

The key idea of probing method is to utilize supervised learning algorithms to probe internal representations in end-to-end neural models, to predict linguistic properties, such as part-of-speech or morphology. Specifically, the method consists of three steps:

1. train an end-to-end model on a complex task, such as machine translation.

2. use the trained model to generate feature representations of different layers.

3. train a classifier using the generated features to make predictions for a relevant auxiliary task, such as POS tagging.

This process is illustrated in Figure 6.1

Formally, let $\mathbf{M}_\theta$ denote a deep neural network with $L$ hidden layers, indexed by $l \in \{1, \ldots, L\}$. Let $\mathbf{h}^{(l)}$ denote the output vector from layer $l$. As usual, $\mathbf{h}^{(0)} = x$ is the input, and $\mathbf{h}^{(L)}$ is the output of the neural network. Denote $\theta = \{\theta_l : l = 1, \ldots, L\}$ as the set of parameters in the network $\mathbf{M}_\theta$, where $\theta_l$ assembles the parameters in layer $l$. Define a separate classifier $g_\phi(\cdot)$ that takes a internal representation $\mathbf{h}^{(l)}$ of $\mathbf{M}_\theta$ as input and maps it to an output

label $z$. After the first step that $\mathbf{M}_\theta$ is trained on the complex task such as machine translation, the parameters $\theta$ are fixed. At the second step, $\mathbf{M}_\theta$ generates internal feature representations like $\mathbf{h}^{(l)}$. At the last step, $g_\phi(\cdot)$ is trained on examples $\{\mathbf{h}^{(l)}, z\}$ to optimize $\phi$. Crucially, at this step $\theta$ is fixed in order to maintain the original generated representation $\mathbf{h}^{(l)}$. More detailed description of probing method is in Belinkov (2018).

## 6.3   Related Work of Probing

Probing classifiers are the most common approach for associating neural network representations with linguistic properties (see Belinkov et al. (2019) for a survey). In their pioneering work, Ettinger et al. (2016) and Shi et al. (2016) investigated intermediate layers of deep neural models in NLP and Alain and Bengio (2016) in computer vision. Both of them used linear classifiers as their probes. Subsequently, extensive applications of probing method have been explored. For the feature representations, probing analysis has been conducted on word embeddings (Köhn, 2015; Qian et al., 2016b) sentence embeddings (Adi et al., 2017; Ganesh et al., 2017; Conneau et al., 2018), and RNN hidden states (Qian et al., 2016a; Wu and King, 2016; Wang et al., 2017). Liu et al. (2019a) and Tenney et al. (2019) focused on contextual word representations, evaluating CoVe (McCann et al., 2017), ELMo (Peters et al., 2018) and BERT (Devlin et al., 2019) on a variety of linguistic tasks. Recently studies carried a more fine-grained neuron-level analysis for neural machine translation and language modeling (Bau et al., 2018; Dalvi et al., 2019; Lakretz et al., 2019). For the language properties, analysis is performed mainly on morphology (Qian et al., 2016b; Vylomova et al., 2017; Belinkov et al., 2017a; Dalvi et al., 2017), syntax (Köhn, 2015; Tran et al., 2018; Conneau et al., 2018; Smith et al., 2018) and semantics (Qian et al., 2016b; Belinkov et al., 2017b). These research attempted to answer the natural question: *what linguistic information is captured in the internal representations*, with the basic assumption that the performance of the commonly used probes such as linear classifiers and multiple layer perceptron (MLP) reflects the quality of representations (Belinkov, 2018; Liu et al., 2019a).

Despite widespread adoption of probes, recent studies illustrated that differences in the accuracy of these commonly used probes accurately fail to adequately reflect the differences in

quantity of the information encoded in the representations. For example, they do not substantially favor pretrained representations over randomly initialized ones (Zhang and Bowman, 2018). Analogously, their accuracy can be similar when probing genuine linguistic labels and probing for random synthetic control tasks (Hewitt and Manning, 2019). To see differences in the accuracy with respect to these random baselines, previous work has examined how the choice of probing tasks and models, and comparing baselines affect the probing conclusion. Some studies (Belinkov et al., 2017a; Tenney et al., 2019; Hewitt and Manning, 2019) used non-contextual word embeddings or models with random weights as baselines. Zhang and Bowman (2018) presented experiments for understanding the roles probe training sample size have on linguistic task accuracy. Hewitt and Liang (2019) designed control tasks to explore the relationship between representations, probes and task accuracies. Saphra and Lopez (2019) showed that diagnostic classifiers are not suitable for understanding learning dynamics.

In addition to task performance of linguistic properties, Yogatama et al. (2019) used learning curves to evaluate how quickly a model learns a new task. Talmor et al. (2019) explored whether the performance of a language model on a task should be attributed to the pre-trained representations or to the process of fine-tuning on the task data.

Information-theoretic view of measuring relations between representations and labels is an alternative to the standard probing. Belinkov (2018) first presented the association of the probabilistic probing classifiers with mutual information between internal representation $\mathbf{h}^{(l)}$ and label $z$. Voita et al. (2019) attempted to explain how representations in the Transformer (Vaswani et al., 2017) evolve between layers under different training objectives. Voita and Titov (2020) proposed information-theoretic probing with minimum description length to evaluate the amount of effort needed to achieve the quality given representations.

## 6.4 Revisit Probing Performance in A View of Prediction

### 6.4.1 Capacity vs. Accessibility

In this section, we revisit this problem by taking a prediction view of the probing accuracy on reflecting the quality of linguistic properties given representations. We argue that we cannot

meaningfully compare the linguistic properties of internal representations of neural networks using only linguistic task accuracy, since probing classifiers with different expressiveness may obtain inconsistent observations. This leads us to re-think about how to interpret performance of probing tasks. on one hand, as long as one representation is a lossless mapping of the other one, a sufficiently expressive probe with sufficient amount of training data can obtain the same accuracy on top of them. However, it fails to take into account differences in ease of memorization between them. On the other hand, if the probe is too weak, then the representations may contain information that cannot be extracted by the classifier, leading to negative results on the quality of $M_\theta$. Specifically, we propose to analyze representations with two separate metrics that are associated with the expressiveness of the probe families:

**Capacity:** *the best prediction performance of linguistic task with arbitrary probes.*

**Accessibility:** *the best prediction performance of linguistic task with linear probes.*

Formally, let $\mathcal{G}$ denotes the set of all classifiers and $\mathcal{G}_{linear}$ is the classifiers of linear model family. $R(\cdot, \cdot)$ is a metric to evaluate the prediction performance of the linguistic task. Then we define

$$
\begin{aligned}
\text{Capacity} &= \max_{g \in \mathcal{G}} \mathrm{E}_{P(\mathbf{h},z)}\left[R(g(\mathbf{h}), z)\right] \\
\text{Accessibility} &= \max_{g \in \mathcal{G}_{linear}} \mathrm{E}_{P(\mathbf{h},z)}\left[R(g(\mathbf{h}), z)\right]
\end{aligned}
\tag{6.1}
$$

As two separate metrics, *capacity* intuitively measures how much knowledge associated with the linguistic property has been encoded in the representation, while *accessibility* measures how easily the encoded knowledge can be detected by a linear probe.

Previous studies (Zhang and Bowman, 2018; Hewitt and Manning, 2019; Voita and Titov, 2020) have discussed the problem on the association between probing performance and linguistic property given representations in different aspects and explored the solutions (see Section 6.3). For example, Belinkov (2018) discussed the potential limitations of the standard probing methodology. The standard probing approach relies on the assumption that the performance of the probing classifier $g$ reflects the quality of the end-to-end neural model $M_\theta$. One potential concern is that the classifier is either too weak or too strong. If it is too weak, then the representations may contain information that cannot be extracted by the classifier, leading to negative results on the quality of $M_\theta$. If the classifier is too strong, it may be able to discover patterns that $M_\theta$ cannot

utilize. Different from it, capacity measures how much knowledge has been encoded, not caring whether this knowledge is utilizable by the neural networks. The probe selectivity associated with control tasks (Hewitt and Liang, 2019), which is a metric of the probe, gives us indirect intuition for the ease of information memorization across different representations. Different from it, our proposed two metrics provide direct insight into how much information has been encoded and how easily the information can be accessed by ML models. Voita and Titov (2020) proposed to use minimum description length to evaluate the "amount of effort" needed to achieve the quality given representations.

## 6.4.2   Probing POS Information in Neural Dependency Parsers

To verify the importance of associating probing accuracy with classifier expressiveness and the necessity of the introduced two metrics, we use probing methods to analyze the part-of-speech (POS) information encoded in the internal representation of three state-of-the-art neural dependency parsing models — Deep Biaffine Parser, NeuroMST parser and Stack-Pointer Parser, . In all experiments, we use SVM-RBF as the probe to approximate capacity and linear logistic regression for accessibility.

We first explore the consistency of using probes with different expressiveness: *Do probing classifiers with different expressiveness lead to consistent observations?* Then we conduct a battery of experiments to investigate the POS information in neural dependency parsers by answering the following questions:

- What is the division of labor between word and character embeddings?

- Which parts of the neural architecture capture POS information?

- Does the addition of accurate POS information impact the learned representations in terms of POS information?

- What impact does the choice of parsing algorithms (graph-based vs. transition-based) have on the learned representations?

All experiments are performed on PTB with the same settings in Ma et al. (2018).

We select dependency parsing as the test bed primarily because of the accessibility to annotated

Figure 6.2: Architecture of the BLSTM-CNN encoder in the three neural parsers. The "char" embedding is the output representation from the CNN layer.

data across a broad spectrum of languages (Nivre et al., 2018) and the sufficiently sophisticated performance of state-of-the-art dependency parsers (Dozat and Manning, 2017; Ma and Hovy, 2017; Ma et al., 2018). we conduct probing experiments on three neural dependency parsing models — Deep Biaffine Parser (Dozat and Manning, 2017), NeuroMST parser (Ma and Hovy, 2017) (in Chapter 4 and Stack-Pointer Parser (Ma et al., 2018) (in Chapter 5). All the three parsing models are implemented with the bi-directional LSTM-CNN architecture (BLSTM-CNN) (Chiu and Nichols, 2016; Ma and Hovy, 2016) as encoder to incorporate both word-level and character-level information. For Deep Biaffine parser we use the re-implemented version in Ma et al. (2018). The BLSTM-CNN encoder consists of three bi-directional LSTM layers (see the architecture depicted in Figure 6.2).

**Motivation: a case study on the effect of POS embedding**    From the results in Section 5.3.4 we observed that end-to-end neural parsers without POS information can obtain even better performance than parsers using imperfect predicted POS tags. To verify this, we performed similar experiments on English Penn Treebank (PTB) (Marcus et al., 1993) using the three parsing models. For each parsing model, we evaluate the parser with gold-standard and predicted POS

Figure 6.3: UAS on the test data of PTB with different versions of POS tags.

tags, and without POS tags, respectively. For predicted POS tags, we used a pretrained POS tagger with 97.3% accuracy.

Figure 6.3 lists the unlabeled attachment score (UAS) of the three parsers with different versions of POS tags on PTB. Consistent with the observation in Ma et al. (2018), the parsers with gold-standard POS tags achieve the best UAS, while the parsers with predicted (imperfect) POS tags perform slightly worse than the parsers that do not use explicit POS tags at all. This observation raises a question: how much POS information is captured implicitly (if anything) by neural parsing models and where this information is stored, which motivates us to investigate the information encoded in the internal representations of these parsers.

### 6.4.3 Word and Character Embeddings

We evaluate the POS tagging performance with three probing classifiers — linear logistic regression (Linear), multiple layer perceptron with one hidden layer (MLP) and Support Vector Machine with RBF kernel (SVM-RBF) — on the word and character embeddings. For comprehensive analysis, we also evaluate the concatenations word and character embeddings (W⊕C).

Table 6.1 shows the probing results with three classifiers, together with MFT, which is the most frequent tag baseline on word types. For out-of-vocabulary (OOV) words, we assign them the most common POS tags. For each classifier, we present the accuracy of both the POS tagging and its control task (Hewitt and Liang, 2019), which associates word types with random outputs to complement linguistic tasks, together with the corresponding selectivity which is the gap between

|  | Linear | | | MLP | | | SVM-RBF | | | MFT | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Layers** | POS | Ctl. | Sel. | POS | Ctl. | Sel. | POS | Ctl. | Sel. | POS | Ctl | Sel. |
| Word | 91.9 | 58.0 | 33.9 | 92.0 | 65.1 | 26.9 | 92.8 | 94.0 | -1.2 | 92.6 | 96.0 | -3.4 |
| Char | 90.3 | 58.3 | 32.0 | 90.5 | 70.6 | 19.9 | 93.5 | 96.1 | -2.6 | – | – | – |
| W⊕C | 93.4 | 65.9 | 27.5 | 93.5 | 74.8 | 18.7 | 93.7 | 96.2 | -2.5 | – | – | – |

Table 6.1: Performance of three probing classifiers on word and character embeddings, and their concatenation of a DeepBiaf parser trained on PTB without POS tags as input. MFT is the most frequent tag baseline. Clt. is the perfomance of the POS tagging control task and Sel. is the corresponding selectivity.

the accuracies of the linguistic task and its control task. Note that MFT can be regarded as an upper bound for non-contextual representations, such as word and character embeddings, for both POS tagging and its control task. The observations of Linear and MLP classifiers on the three representations are similar, with the difference on the selectivity of the two probes. According to the performance of Linear and MLP classifiers, the word representation encodes more POS information than the character one, while the concatenated representation obtains much better accuracy, even better than the MFT upper bound. It demonstrates that character embedding encodes more POS information than the word embedding. But the POS information in the word embedding is easier to be detected by probes — a simple linear probing classifier obtains higher accuracy on top of the word embedding. If we only consider those results from Linear and MLP classifiers, we may conclude that the POS information encoded in word and character embeddings are complementary to each other. However, the performance of SVM-RBF classifier provides significantly different observations: (i) the character embedding achieves better accuracy than the word embedding (93.5% vs. 92.8%); (ii) the concatenated representation achieves similar accuracy with the character embedding (93.7% vs. 93.5%), leading to the conclusion that POS information in the word embedding is mostly covered by the character embedding.

### 6.4.4 Effects of POS Embedding as Input

Table 6.2 illustrates the accessibility and capacity on representations from different layers of DeepBiaf parsers trained w./wo. POS tags as input. We observe that POS embedding has

|  | Word | | Char | | W⊕C | | LSTM 1 | | LSTM 2 | | LSTM 3 | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **POS** | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. |
| WO. | 91.9 | 92.8 | **90.3** | 93.5 | 93.4 | 93.7 | 97.5 | 97.7 | 97.4 | 97.8 | 95.7 | 96.8 |
| W. | 91.8 | 92.8 | **89.5** | 93.5 | 93.5 | 93.7 | 100.0 | 100.0 | 99.7 | 99.7 | 97.3 | 98.4 |

Table 6.2: Accessibility (Acces.) and Capacity (Cap.) on representations from different layers of DeepBiaf parsers trained on PTB w./wo. POS tags as input. Qualities on character embeddings are highlighted because of the significant difference.

|  | Word | | Char | | W⊕C | | LSTM 1 | | LSTM 2 | | LSTM 3 | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **Parsers** | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. | Acces. | Cap. |
| **DeepBiaf** | 91.9 | 92.8 | 90.3 | 93.5 | 93.4 | 93.7 | 97.5 | 97.7 | 97.4 | 97.8 | 95.7 | 96.8 |
| **NeuroMST** | 91.9 | 92.8 | 90.1 | 93.5 | 93.4 | 93.7 | 97.5 | 97.6 | 97.4 | 97.8 | 96.0 | 96.9 |
| **StackPtr** | 91.8 | 92.9 | 90.1 | 93.5 | 93.4 | 93.7 | 97.5 | 97.7 | 97.1 | 97.7 | **94.0** | **96.1** |

Table 6.3: Accessibility (Acces.) and Capacity (Cap.) on representations from the layers of three parsers. Results that are significantly different from that of the same layers of other parsers are highlighted.

no significant impact on word or character embeddings, except the accessibility of character embedding — without using POS embedding the accessibility of character embedding on POS information improves.

We also observe that even without using POS embedding, the LSTM layers of the encoder capture a large amount of POS information. Both the accessibility and capacity of the first and second LSTM layers are better than the accuracy (97.3%) of the predicted POS tags, explaining why the parsers without using POS tags achieved better performance than those using predicted POS tags. In fact, the POS tagging accuracy of 97.8% (the capacity of the second LSTM layer) is competitive or even better than the state-of-the-art neural POS tagging systems (Ma and Hovy, 2016; Cui and Zhang, 2019). It verifies that learning a more complex task (dependency parsing) benefits the simpler upstream tasks (POS tagging).

Another interesting observation is that, for both parsers trained w./wo. POS embedding, the POS information captured in the top LSTM layer (LSTM 3) declines in both accessibility and

capacity comparing with the first two layers. One possible explanation for the discrepancy is that the top LSTM layer, which is directly used for the parsing, pays more attention to parsing task.

### 6.4.5   Effects of Parsing Algorithms

Table 6.3 presents the results on representations from layers of three dependency parsers. For the two graph-based parsers with different training objectives (edge-factored loss of DeepBiaf vs. global structured loss of NeuroMST), there is no significant difference on either the accessibility or the capacity of the representations.

When comparing the transition-based StackPtr with the two graph-based parsers (DeepBiaf and NeuroMST), we observe that for non-contextual representations (word and character embeddings) and the first two LSTM layers, the capacity and accessibility are similar. But when the depth increases, both the capacity and accessibility of the POS information in the transition-based StackPtr parser declines more rapidly than the two graph-based parsers, demonstrating that different decoding algorithms indeed impact the output representations of the encoder.

## 6.5   Discussions

In this chapter, we probed the internal representations of three neural dependency parsers. Through analyzing the results of probes in different expressive families, we found that these probes obtained significant different accuracies. Based on this, we proposed two complementary metrics: Capacity and Accessibility, to enhance the interpretation of results from probing methods.

In the next chapter, we will introduce how to use probes to investigate learning properties of deep neural networks, rather than analyzing what linguistic information is captured in their internal layers.

# Chapter 7

# Probing Learning Properties of Neural Networks

As discussed in the previous chapter, the research of probing methods attempts to answer the natural question: what linguistic information is captured in the internal representations of neural networks. Moreover, besides using probes to analyze what linguistic information is captured in neural networks, can we use probes to understand the way deep neural networks memorize and process linguistic information? To answer this question, Shi et al. (2016) conducted probing experiments to explore what types of information are learned with different training objectives. Saphra and Lopez (2019) presented the analysis on learning dynamics of neural language models to compare learned representations across time and models, without the evaluation directly on annotated data. Learning curves have also been used by Yogatama et al. (2019) to evaluate how quickly a model learns a new task, and by Talmor et al. (2019) to understand whether the performance of a language model on a task should be attributed to the pre-trained representations or to the process of fine-tuning on the task data. Despite these previous attempts, this question still remains unclear.

In this chapter, we use probing methods to investigate the learning properties of deep neural networks with dependency parsing as a test bed. By conducting systematic experiments, we illustrate two learning properties of deep neural networks: (i) *laziness* — modules of a neural network will not actively learn information that is already learned by other modules (Section 7.1);

Figure 7.1: An example for the illustration of our designed experiments to examine the laziness of deep neural networks.

(ii) *targetedness* — information, if unnecessary for the target task, will be filtered out from the internal representations (Section 7.2).

# 7.1 Laziness: Information Propagation through Deep Neural Networks

In this section, we investigate one learning property of neural networks by analyzing the information propagation during neural network training. Results from carefully designed probing experiments illustrate that deep neural networks exhibit the learning property of *laziness* — information, if already encoded in some components of a neural network, will not be propagated to the network's other components.

## 7.1.1 Laziness vs. Redundancy

In this work, laziness, or the opposite redundancy, is defined on the pattern of how neural networks distribute information across their components — do they store one kind of information/knowledge in a single module or multiple ones?

Figure 7.2: The heatmap of capacity and accessibility of DeepBiaf parsers with POS tags as input to different layers. pos1, pos2 and pos3 refer to the parser with gold-standard POS tags as input to the 1st, 2nd and 3rd LSTM layers. none indicates the parser without POS tags as input.

As discussed in the pioneering of Dropout (Srivastava et al., 2014a), redundancy is commonly a desired property and one of the primary motivations of the dropout approach. Redundancy is closely related to robustness, because a system that stores redundant information has better chance to work robustly under non-stationary conditions. Laziness, on the other hand, is probably the way that requires minimal amount of efforts to learn knowledge or to fullfil target tasks. Dividing the target task into smaller individual sub-tasks, though might not be robust, may be the most efficient way to complete the ultimate goal.

### 7.1.2 Experiment Design

The key idea of the experiments is to examine the learning behaviors of neural parsers by placing the gold-standard POS tags as input to different encoder layers. Suppose that the gold-standard POS tags are fed as the input of the third LSTM layer (see Figure 7.1 for illustration). Then, during the forward pass to compute the parsing objective, the gold-standard POS information will be passed to the third LSTM layer and the layers on top of it, but not the layers below it. During the backward pass in training, however, the gradients that might carry gold-standard POS information will be back-propagated to the entire neural network. By probing the layers of this neural parser,

we can examine if the gold-standard POS information has been stored in the layers below the one that gold-standard POS tags are fed into or not. If the gold-standard POS information has been propagated to the lower layers, as illustrated in the second case in the right side of Figure 7.1, the POS accuracy of probes should be significantly better than the one without using POS tags, showing the redundancy. Otherwise, it demonstrates the laziness.

In this experiment, we train DeepBiaf parsers on PTB with gold-standard POS tags as input to different LSTM layers. Dropout (Srivastava et al., 2014a) is applied everywhere (see Ma et al. (2018) for details). Figure 7.2 shows the heatmap of the capacity and accessibility. From the heatmap, we can see clearly that the probing accuracies in the layers below the one that POS tags are fed into are significantly worse than that of the above layers. It implies that the POS information has not been propagated into those lower layers through back propagation, verifying the laziness of neural networks. Importantly, dropout, which is introduced to enhance the redundancy of neural networks, does not as expected prevent the laziness learning of the neural parsers in our experiments. What the real impact dropout has on neural network training remains an open problem and might be an interesting direction for future work (Ma et al., 2017).

## 7.2 Targetedness: Semantic Information in Syntactic Parsers

In this section, we explore another learning property of neural networks, *targetedness*, by asking the question how deep neural networks process information that is unnecessary for their ultimate tasks — will this information be retained in their internal layers or filtered out?

### 7.2.1 Lexical Semantic Tagging

To answer the above question, we conduct experiments to probe neural dependency parsers with lexical semantic tagging (SEM) (Bjerva et al., 2016) as the linguistic probing task. SEM is a sequence labeling task: given a sentence, the goal is to assign to each word a tag representing a semantic class. Figure 7.3 gives an example of SEM. Linguistically, some lexical semantic information is unnecessary for syntactic parsing. For instance, w.r.t dependency parsing, the semantic tag PER (for person) of the token "Tom" provides no more useful information than its

| SEM | PER | NOW | NOT | EXS | QUV | CON | . |
|---|---|---|---|---|---|---|---|
| POS | NNP | VBZ | RB | VB | JJ | NN | . |
| | Tom | does | not | drink | much | beer | . |

Figure 7.3: An example of lexical semantic tagging.

POS tag NNP (for proper noun). With SEM as the probing task, we want to examine if these unnecessary lexical semantic information will be filtered out from dependency parsing models trained with only syntactic supervision.

For the annotated data of lexical semantic tagging, we use the Groningen Parallel Meaning Bank (PMB) (Abzianidze et al., 2017) which includes 66 fine-grained semantic tags grouped in 13 coarse categories. The experiments are conducted on the silver part of the dataset — we randomly split the data into training, validation and test sets with the proportion of $[8, 1, 1]$.

## 7.2.2 Baselines

In order to analyze the lexical semantic information memorized in the neural networks, we need reasonable baselines for comparison. In Belinkov (2018), the authors proposed two baselines: (i) assigning to each word the most frequent tag (MFT) according to the training data, with the global majority tag for unseen words; (ii) training with unsupervised word embeddings as features for the classifier. However, these two baselines are arguably unsuitable for our experiments, since they are based on non-contextual features. Comparing these two baselines with the internal contextual representations in neural dependency parsers is unfair and cannot lead to reliable conclusions. In other words, even if the probes on the internal representations achieve significant better accuracy than the two baselines, we cannot conclude that the neural dependency parser has learned lexical semantics. The reason is that neural parsers are able to learn POS information (as shown in Section 6.4.2), which is highly correlated with lexical semantic information (assigning to each word the most frequent tag based on its POS tag obtains 75.2% accuracy). Thus, to claim that the neural parsers cant capture lexical semantic information, we need to proof that the representations capture more information than that provided by POS tags.

In this paper, we propose two new baselines: (i) assigning to each word the most frequent

|  |  | SEM Accuracy |  |
| --- | --- | --- | --- |
| MFT (word) |  | 85.0 |  |
| MFT (word + POS) |  | 89.9 |  |
| BF-Linear |  | **93.4** |  |
| BLSTM-CNN-CRF |  | **94.9** |  |
|  |  | Acces. | Cap. |
|  | Word | 82.3 | 83.8 |
|  | Char | 80.3 | 85.2 |
| **Layers** | W⊕C | 85.4 | 85.3 |
|  | LSTM 1 | 91.5 | **91.7** |
|  | LSTM 2 | 90.1 | 91.6 |
|  | LSTM 3 | 86.8 | 89.5 |

Table 7.1: Accessibility and Capacity of lexical semantic tagging on representations from layers of DeepBiaf dependency parsers, together with MFT baselines, the linear classifier on binary features (BF-Linear) and the upper bound accuracy of the sequence labeling model.

SEM tag based on the combination of its word type and POS tag; (ii) a linear classifiers trained on binary features of word type of each token and context of POS tags in a small neighborhood (window of 3). The POS tags are automatically labeled with Standford POS Tagger[1]. Note that neither of these two baselines utilize contextual information beyond POS tags. We also include an upper bound of training a BLSTM-CNN-CRF sequence labeling model (Ma and Hovy, 2016) for SEM tagging.

### 7.2.3 Experimental Results

The experiments are conducted on DeepBiaf parsers trained on PTB without using POS tags as input. Table 7.1 summarizes the results of training probes to predict SEM tags using representations generated by different encoding layers of the DeepBiaf parser. Comparing representations from

[1]https://nlp.stanford.edu/software/tagger.shtml

LSTM layers 1 through 3, SEM tagging accuracy peaks at layer 1 and does not improve at higher layers. The best capacity score from LSTM layer 1 is 91.7%, slightly worse than the baseline of linear classifier with binary features, far below the upper bound from the sequence labeling model. It indicates that the internal representations do not learn more contextual information that benefits SEM prediction than POS tags.

### 7.2.4 Discussion

Extensive studies investigated the requested semantic information in syntactic parsing, such as selectional restrictions (Katz and Fodor, 1963; Wilks, 1975; Jurafsky and Martin, 2000; Asher, 2014). Lexical semantic categories, in some cases, are indeed relevant for selectional restrictions and in turn can be used for disambiguation in syntactic parsing. One possible reason that neural dependency parsers are not able to learn relevant semantic information is that there are few such cases in the training data. Another possible reason might be due to the SEM data. The PMB data used in our experiments do not support analysis on specific case studies, but only an average accuracy on all the semantic categories. Fine-grained investigation on specific categories of semantic information learned by neural parsers is an interesting direction for future work.

One previous work that is closely related to the property of targetedness is Shi et al. (2016), which demonstrated that when training with different objectives, neural networks capture different types of information. This observation is relevant to targetedness, but not exactly the same. This dissertation first (to our best knowledge) explicitly demonstrates that neural networkss filter out unnecessary information, by specifically control the modifications to neural architectures.

## 7.3 Conclusion

In this chapter, using the the two metrics proposed in Section 6.4.1, we have conducted experiments to investigate the learning properties of neural networks. Experimental results illustrate two learning behaviors of neural networks: (i) laziness — modules of a neural network will not actively learn information that is already learned by other modules; (ii) targetedness — information that is unnecessary for the ultimate objective will be filtered out.

# Chapter 8

# Conclusions and Future Work

## 8.1 Conclusion

The work presented in the first part of this dissertation explores how the deep representation learning approach, particularly the end-to-end learning paradigm, can be applied to linguistic structured predictions to effectively improve model performance and entirely get rid of feature engineering in traditional feature-based models. In PART I of this dissertation, we gave readers a thorough overview of neural networks for linguistic structured prediction: three neural models for sequence and tree structured prediction tasks. These neural models are all base on the proposed BLSTM-CNN architecture for sentence encoding. BLSTM-CNNs provide sentence representations that are applicable across different structured prediction tasks, while supporting the kind of end-to-end learning, saving us from hand-crafted feature engineering (Chapter 2). By stacking different structured decoding layers on top of BLSTM-CNNs, we proposed deep neural models for linguistic structured prediction tasks including sequence labeling (BLSTM-CNNs-CRF in Chapter 3), graph-based dependency parsing (NeuraoMST Parser in Chapter 4) and transition-based dependency parsing (Stack-Pointer Parser in Chapter 5). Experimental results demonstrate that these BLSTM-CNNs based models have achieved significant improvements over traditional feature-based models, and state-of-the-art or comparable performance across different languages and corpus.

In PART II of this dissertation, we focused the interpretability of neural dependency parsing

models . We first revisit the problem of standard probing method on the reflection of probing performance on linguistic properties. Then, we explicitly proposed two separate metrics, capacity and accessibility, by taking a prediction view of probing accuracy. Experimental results on probing how part-of-speech information in neural dependency parsers illustrate the necessity of the two metrics (Chapter 6). We further we conducted systematic experiments to illustrate two learning properties of deep neural networks: (i) *laziness* – modules of a neural network will not actively learn information that is already learned by other modules; (ii) *targetedness* – information, if unnecessary for the end task, will be filtered out from the internal representations. These two properties help us better understand how deep neural structured prediction models memorize and process linguistic information (Chapter 7).

## 8.2   Future Work

This thesis opens up several questions for future research.

**Encoding Structured Dependencies in Representations: No Structured Algorithms in Structured Predictions**

Although the neural networks proposed in this thesis obtained outstanding performance for a wide range of tasks and languages, they still suffer some problems: (i) the design and combination of the structured output layers with the end-to-end neural representation encoders is not easy. Every step of the structured algorithms must be ensured differentiable so that the gradients can be back-propagated to the entire network to achieve end-to-end learning; (ii) most of the decoding algorithms are inefficient in practice because they are not parallelizable.

To fundamentally solve these problems, one potential direction is to encode the underlying dependencies of the structured outputs into intermediate representations to get rid of structured training and decoding algorithms. Ma et al. (2019) proposed to a non-autoregressive sequence generation model to avoid sequential decoding algorithm. On direction of future research is to extend this framework to general structured prediction tasks.

**Investigating Information Encoding Schema in Representations**

Another interesting direction of future research is to investigate the information encoding schema

of different neural architectures. From the observations in our probing experiments, different representations, even encoding the same knowledge, may encode them in different format. The different information encoding schema is highly correlated with the learning properties of the neural architecture. Therefore, figuring out the relation between encoding schema and the neural architectures is useful for us to better understand different neural architectures.

**Inductive Bias from Learning Properties of Different Neural Architectures**

Finally, on promising direction of future research is to explore useful inductive bias from investigating the learning properties of different neural architectures. By investigating the learning properties, such as the two learning properties of we investigated in Chapter 7, we can conclude architectural inductive biases and attempt to apply them to advanced representation learning, such as disentanglement, to get rid of explicit supervision.

# Appendix A

# Dependency Parsing Experiments

## A.1 Hyper-parameters

### A.1.1 NeuroMST Parser

Table A.1 summarizes the chosen hyper-parameters for NeuroMST parser. We tune the hyper-parameters on the development sets by random search. Due to time constrains it is infeasible to do a random search across the full hyper-parameter space. Thus, we use the same hyper-parameters across the models on different treebanks and languages. It also demonstrates the robustness of our parsing model. Note that we use 2-layer BLSTM followed with 1-layer MLP. We set the state size of LSTM to $256$ and the dimension of MLP to $100$. Tuning these two parameters did not significantly impact the performance of our model.

| Layer | Hyper-parameter | Value |
|---|---|---|
| CNN | window size | 3 |
| | number of filters | 50 |
| LSTM | number of layers | 2 |
| | state size | 256 |
| | initial state | 0.0 |
| | peepholes | Hadamard |
| MLP | number of layers | 1 |
| | dimension | 100 |
| Dropout | embeddings | 0.15 |
| | LSTM hidden states | 0.25 |
| | LSTM layers | 0.33 |
| Learning | optimizer | Adam |
| | initial learning rate | 0.002 |
| | decay rate | 0.5 |
| | gradient clipping | 5.0 |

Table A.1: Hyper-parameters for NeuroMST parser.

| Layer | Hyper-parameter | Value |
|---|---|---|
| CNN | window size | 3 |
| | number of filters | 50 |
| LSTM | encoder layers | 3 |
| | encoder size | 512 |
| | decoder layers | 1 |
| | decoder size | 512 |
| MLP | arc MLP size | 512 |
| | label MLP size | 128 |
| Dropout | embeddings | 0.33 |
| | LSTM hidden states | 0.33 |
| | LSTM layers | 0.33 |
| Learning | optimizer | Adam |
| | initial learning rate | 0.001 |
| | $(\beta_1, \beta_2)$ | (0.9, 0.9) |
| | decay rate | 0.75 |
| | gradient clipping | 5.0 |

Table A.2: Hyper-parameters for StackPtr parser.

## A.1.2 StackPtr Parser

Table A.2 summarizes the chosen hyper-parameters used for all the experiments in this paper. Some parameters are chosen directly or similarly from those reported in Dozat and Manning (2017). We use the same hyper-parameters across the models on different treebanks and languages, due to time constraints.

## A.2   UD Treebanks

The statistics of these corpora are provided in Table A.3.

| | Corpora | | #Sent | #Token (w.o punct) |
|---|---|---|---|---|
| Bulgarian | BTB | Training | 8,907 | 124,336 (106,813) |
| | | Dev | 1,115 | 16,089 (13,822) |
| | | Test | 1,116 | 15,724 (13,456) |
| Catalan | AnCora | Training | 13,123 | 417,587 (371,981) |
| | | Dev | 1,709 | 56,482 (50,452) |
| | | Test | 1,846 | 57,738 (51,324) |
| Czech | PDT, CAC | Training | 102,993 | 1,806,230 (1,542,805) |
| | CLTT | Dev | 11,311 | 191,679 (163,387) |
| | FicTree | Test | 12,203 | 205,597 (174,771) |
| Dutch | Alpino | Training | 18,310 | 267,289 (234,104) |
| | LassySmall | Dev | 1,518 | 22,091 (19,042) |
| | | Test | 1,396 | 21,126 (18,310) |
| English | EWT | Training | 12,543 | 204,585 (180,308) |
| | | Dev | 2,002 | 25,148 (21,998) |
| | | Test | 2,077 | 25,096 (21,898) |
| French | GSD | Training | 14,554 | 356,638 (316,780) |
| | | Dev | 1,478 | 35,768 (31,896) |
| | | Test | 416 | 10,020 (8,795) |
| German | GSD | Training | 13,841 | 263,536 (229,204) |
| | | Dev | 799 | 12,348 (10,727) |
| | | Test | 977 | 16,268 (13,929) |
| Italian | ISDT | Training | 12,838 | 270,703 (239,836) |
| | | Dev | 564 | 11,908 (10,490) |
| | | Test | 482 | 10,417 (9,237) |
| Norwegian | Bokmaal | Training | 29,870 | 48,9217 (43,2597) |
| | Nynorsk | Dev | 4,300 | 67,619 (59,784) |
| | | Test | 3,450 | 54,739 (48,588) |
| Romanian | RRT | Training | 8,043 | 185,113 (161,429) |
| | | Dev | 752 | 17,074 (14,851) |
| | | Test | 729 | 16,324 (14,241) |
| Russian | SynTagRus | Training | 48,814 | 870,034 (711,184) |
| | | Dev | 6,584 | 118,426 (95,676) |
| | | Test | 6,491 | 117,276 (95,745) |
| Spanish | GSD | Training | 28,492 | 827,053 (730,062) |
| | AnCora | Dev | 4,300 | 89,487 (78,951) |
| | | Test | 2,174 | 64,617 (56,973) |

Table A.3: Corpora statistics of UD Treebanks for 12 languages. *#Sent* and *#Token* refer to the number of sentences and the number of words (w./w.o punctuations) in each data set, respectively.

## A.3   Detailed Results of Stack-Pointer Parser

Table A.4 illustrates the details of the experimental results. For each STACKPRT parsing model, we ran experiments with decoding beam size equals to 1, 5, and 10. For each experiment, we report the mean values with corresponding standard deviations over 5 runs.

| Model | beam | English | | | | | | | |
| | | Dev | | | | Test | | | |
| | | UAS | LAS | UCM | LCM | UAS | LAS | UCM | LCM |
|---|---|---|---|---|---|---|---|---|---|
| BiAF | – | 95.73±0.04 | **93.97±0.06** | 60.58±0.77 | 47.47±0.63 | 95.84±0.06 | **94.21±0.04** | 59.49±0.23 | 49.07±0.34 |
| Basic | 1 | 95.71±0.02 | 93.88±0.03 | 62.33±0.33 | 47.75±0.32 | 95.71±0.06 | 94.07±0.06 | 60.91±0.35 | 49.54±0.48 |
| | 5 | 95.71±0.04 | 93.88±0.05 | 62.40±0.45 | 47.80±0.44 | 95.76±0.11 | 94.12±0.11 | 61.09±0.43 | 49.67±0.41 |
| | 10 | 95.72±0.03 | 93.89±0.04 | **62.40±0.45** | **47.80±0.44** | 95.77±0.11 | 94.12±0.11 | 61.09±0.43 | 49.67±0.41 |
| +gpar | 1 | 95.68±0.04 | 93.82±0.02 | 61.82±0.36 | 47.32±0.14 | 95.73±0.04 | 94.07±0.05 | 60.99±0.46 | 49.83±0.59 |
| | 5 | 95.67±0.01 | 93.83±0.02 | 61.93±0.32 | 47.44±0.20 | 95.76±0.06 | 94.11±0.06 | 61.23±0.47 | 50.07±0.59 |
| | 10 | 95.69±0.02 | 93.83±0.02 | 61.95±0.32 | 47.44±0.20 | 95.78±0.05 | 94.12±0.06 | 61.24±0.46 | **50.07±0.59** |
| +sib | 1 | 95.75±0.03 | 93.93±0.04 | 61.93±0.49 | 47.66±0.48 | 95.77±0.15 | 94.11±0.06 | 61.32±0.37 | 49.75±0.29 |
| | 5 | 95.74±0.02 | 93.93±0.05 | 62.16±0.22 | 47.68±0.54 | 95.84±0.09 | 94.17±0.09 | 61.52±0.57 | 49.91±0.76 |
| | 10 | 95.75±0.02 | 93.94±0.06 | 62.17±0.20 | 47.68±0.54 | 95.85±0.10 | 94.18±0.09 | **61.52±0.57** | 49.91±0.76 |
| Full | 1 | 95.63±0.08 | 93.78±0.08 | 61.56±0.63 | 47.12±0.36 | 95.79±0.06 | 94.11±0.06 | 61.02±0.31 | 49.45±0.23 |
| | 5 | 95.75±0.06 | 93.90±0.08 | 62.06±0.42 | 47.43±0.36 | 95.87±0.04 | 94.20±0.03 | 61.43±0.49 | 49.68±0.47 |
| | 10 | **95.75±0.06** | 93.90±0.08 | 62.08±0.39 | 47.43±0.36 | **95.87±0.04** | 94.19±0.04 | 61.43±0.49 | 49.68±0.47 |
| Model | beam | Chinese | | | | | | | |
| | | Dev | | | | Test | | | |
| | | UAS | LAS | UCM | LCM | UAS | LAS | UCM | LCM |
| BiAF | – | 90.20±0.17 | 88.94±0.13 | 43.41±0.83 | 38.42±0.79 | 90.43±0.08 | 89.14±0.09 | 42.92±0.29 | 38.68±0.25 |
| Basic | 1 | 89.76±0.32 | 88.44±0.28 | 45.18±0.80 | 40.13±0.63 | 90.04±0.32 | 88.74±0.40 | 45.00±0.47 | 40.12±0.42 |
| | 5 | 89.97±0.13 | 88.67±0.14 | 45.33±0.58 | 40.25±0.65 | 90.46±0.15 | 89.17±0.18 | 45.41±0.48 | 40.53±0.48 |
| | 10 | 89.97±0.14 | 88.68±0.14 | 45.33±0.58 | 40.25±0.65 | 90.48±0.11 | 89.19±0.15 | 45.44±0.44 | 40.56±0.43 |
| +gpar | 1 | 90.05±0.14 | 88.71±0.16 | 45.63±0.52 | 40.45±0.61 | 90.28±0.10 | 88.96±0.10 | 45.26±0.59 | 40.38±0.43 |
| | 5 | 90.17±0.14 | 88.85±0.13 | 46.03±0.53 | 40.69±0.55 | 90.45±0.15 | 89.14±0.14 | 45.71±0.46 | 40.80±0.26 |
| | 10 | 90.18±0.16 | 88.87±0.14 | **46.05±0.58** | 40.69±0.55 | 90.46±0.16 | 89.16±0.15 | 45.71±0.46 | 40.80±0.26 |
| +sib | 1 | 89.91±0.07 | 88.59±0.10 | 45.50±0.50 | 40.40±0.48 | 90.25±0.10 | 88.94±0.12 | 45.42±0.52 | 40.54±0.69 |
| | 5 | 89.99±0.05 | 88.70±0.09 | 45.55±0.36 | 40.37±0.14 | 90.41±0.07 | 89.12±0.07 | 45.76±0.46 | 40.69±0.52 |
| | 10 | 90.00±0.04 | 88.72±0.09 | 45.58±0.32 | 40.37±0.14 | 90.43±0.09 | 89.15±0.10 | 45.75±0.44 | 40.68±0.50 |
| Full | 1 | 90.21±0.15 | 88.85±0.15 | 45.83±0.52 | 40.54±0.60 | 90.36±0.16 | 89.05±0.15 | 45.60±0.33 | 40.73±0.23 |
| | 5 | 90.23±0.13 | 88.89±0.14 | 46.00±0.54 | **40.75±0.64** | 90.58±0.12 | 89.27±0.11 | **46.20±0.26** | **41.25±0.22** |
| | 10 | **90.29±0.13** | **88.95±0.13** | 46.03±0.54 | **40.75±0.64** | **90.59±0.12** | **89.29±0.11** | 46.20±0.26 | **41.25±0.22** |
| Model | beam | German | | | | | | | |
| | | Dev | | | | Test | | | |
| | | UAS | LAS | UCM | LCM | UAS | LAS | UCM | LCM |
| BiAF | – | **93.60±0.13** | **91.96±0.13** | 58.79±0.25 | 49.59±0.19 | **93.85±0.07** | **92.32±0.06** | 60.60±0.38 | 52.46±0.24 |
| Basic | 1 | 93.35±0.14 | 91.58±0.17 | 59.64±0.78 | 49.75±0.64 | 93.39±0.09 | 91.85±0.09 | 61.08±0.31 | 52.21±0.53 |
| | 5 | 93.49±0.14 | 91.72±0.16 | 59.99±0.69 | 49.82±0.54 | 93.61±0.09 | 92.07±0.08 | 61.38±0.30 | 52.51±0.43 |
| | 10 | 93.48±0.14 | 91.71±0.17 | **60.02±0.69** | 49.84±0.54 | 93.59±0.09 | 92.06±0.08 | 61.38±0.30 | 52.51±0.43 |
| +gpar | 1 | 93.39±0.07 | 91.66±0.13 | 59.59±0.54 | 49.81±0.42 | 93.44±0.07 | 91.91±0.11 | 61.73±0.47 | 52.84±0.48 |
| | 5 | 93.47±0.09 | 91.75±0.10 | 59.81±0.55 | 50.05±0.39 | 93.68±0.04 | 92.16±0.04 | 62.09±0.44 | 53.13±0.42 |
| | 10 | 93.48±0.08 | 91.76±0.09 | 59.89±0.59 | 50.09±0.40 | 93.68±0.05 | 92.16±0.03 | 62.10±0.42 | **53.14±0.4** |
| +sib | 1 | 93.43±0.07 | 91.73±0.08 | 59.68±0.25 | 49.93±0.30 | 93.55±0.07 | 92.00±0.08 | 61.90±0.50 | 52.79±0.22 |
| | 5 | 93.53±0.05 | 91.83±0.07 | 59.95±0.23 | 50.14±0.39 | 93.75±0.09 | 92.20±0.08 | 62.21±0.38 | 53.03±0.18 |
| | 10 | 93.55±0.06 | 91.84±0.07 | 59.96±0.24 | **50.15±0.40** | 93.76±0.09 | 92.21±0.08 | **62.21±0.38** | 53.03±0.18 |
| Full | 1 | 93.33±0.13 | 91.60±0.16 | 59.78±0.32 | 49.79±0.29 | 93.50±0.04 | 91.91±0.11 | 61.80±0.28 | 52.95±0.37 |
| | 5 | 93.42±0.11 | 91.69±0.12 | 59.90±0.27 | 49.94±0.35 | 93.64±0.03 | 92.10±0.06 | 61.89±0.21 | 53.06±0.36 |
| | 10 | 93.40±0.11 | 91.67±0.12 | 59.90±0.27 | 49.94±0.35 | 93.64±0.03 | 92.11±0.05 | 61.89±0.21 | 53.06±0.36 |

Table A.4: Parsing performance of different variations of our model on both the development and

# Bibliography

Lasha Abzianidze, Johannes Bjerva, Kilian Evang, Hessel Haagsma, Rik van Noord, Pierre Lud-
mann, Duc-Duy Nguyen, and Johan Bos. The parallel meaning bank: Towards a multilingual
corpus of translations annotated with compositional meaning representations. In *Proceedings of
the 15th Conference of the European Chapter of the Association for Computational Linguistics:
Volume 2, Short Papers*, volume 2, pages 242–247, 2017.

Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. Fine-grained
analysis of sentence embeddings using auxiliary prediction tasks. In *Proceedings of ICLR-2017
(Volume 1: Long Papers)*, Toulon, France, August 2017.

Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. Polyglot: Distributed word representations
for multilingual nlp. In *Proceedings of CoNLL-2013*, pages 183–192, Sofia, Bulgaria, August
2013.

Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier
probes. In *Proceedings of the 4th International Conference on Learning Representations
(ICLR-2016)*, April 2016.

Rie Kubota Ando and Tong Zhang. A framework for learning predictive structures from multiple
tasks and unlabeled data. *The Journal of Machine Learning Research*, 6:1817–1853, 2005.

Daniel Andor, Chris Alberti, David Weiss, Aliaksei Severyn, Alessandro Presta, Kuzman Ganchev,
Slav Petrov, and Michael Collins. Globally normalized transition-based neural networks. In
*Proceedings of ACL-2016 (Volume 1: Long Papers)*, pages 2442–2452, Berlin, Germany,
August 2016.

Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D. Manning. Leveraging

linguistic structure for open domain information extraction. In *Proceedings of ACL-2015 (Volume 1: Long Papers)*, pages 344–354, Beijing, China, July 2015.

Nicholas Asher. Selectional restrictions, types and categories. *Journal of Applied Logic*, 12(1): 75–87, 2014.

Giuseppe Attardi. Experiments with a multilanguage non-projective dependency parser. In *Proceedings of the Tenth Conference on Computational Natural Language Learning (CoNLL-X)*, pages 166–170, 2006.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR-2015*, 2015.

Miguel Ballesteros, Chris Dyer, and Noah A. Smith. Improved transition-based parsing by modeling characters instead of words with lstms. In *Proceedings of EMNLP-2015*, pages 349–359, Lisbon, Portugal, September 2015.

Miguel Ballesteros, Yoav Goldberg, Chris Dyer, and Noah A. Smith. Training with exploration improves a greedy stack lstm parser. In *Proceedings of EMNLP-2016*, pages 2005–2010, Austin, Texas, November 2016.

Joost Bastings, Ivan Titov, Wilker Aziz, Diego Marcheggiani, and Khalil Simaan. Graph convolutional encoders for syntax-aware neural machine translation. In *Proceedings of EMNLP-2017*, pages 1957–1967, Copenhagen, Denmark, September 2017.

Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass. Identifying and controlling important neurons in neural machine translation. *arXiv preprint arXiv:1811.01157*, 2018.

Yonatan Belinkov. *On internal language representations in deep learning: an analysis of machine translation and speech recognition*. PhD thesis, Massachusetts Institute of Technology, 2018.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. What do neural machine translation models learn about morphology? In *Proceedings of ACL-2017 (Volume 1: Long Papers)*, pages 861–872, 2017a.

Yonatan Belinkov, Lluís Màrquez, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass.

Evaluating layers of representation in neural machine translation on part-of-speech and semantic tagging tasks. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1–10, 2017b.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. On the linguistic representational power of neural machine translation models. *Computational Linguistics*, pages 1–57, 2019.

Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on*, 5(2):157–166, 1994.

Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.

Taylor Berg-Kirkpatrick, Alexandre Bouchard-Côté, John DeNero, and Dan Klein. Painless unsupervised learning with features. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 582–590. Association for Computational Linguistics, 2010.

James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a cpu and gpu math expression compiler. In *Proceedings of the Python for scientific computing conference (SciPy)*, volume 4, page 3. Austin, TX, 2010.

Daniel M Bikel, Scott Miller, Richard Schwartz, and Ralph Weischedel. Nymble: a high-performance learning name-finder. In *Proceedings of the fifth conference on Applied natural language processing*, pages 194–201. Association for Computational Linguistics, 1997.

Daniel M Bikel, Richard Schwartz, and Ralph M Weischedel. An algorithm that learns what's in a name. *Machine learning*, 34(1-3):211–231, 1999.

Johannes Bjerva, Barbara Plank, and Johan Bos. Semantic tagging with deep residual networks. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3531–3541, 2016.

Bernd Bohnet and Joakim Nivre. A transition-based system for joint part-of-speech tagging and

labeled non-projective dependency parsing. In *Proceedings of EMNLP-2012*, pages 1455–1465, Jeju Island, Korea, July 2012.

Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152, 1992.

Eric Brill. A simple rule-based part of speech tagger. In *Proceedings of the third conference on Applied natural language processing*, pages 152–155. Association for Computational Linguistics, 1992.

Sabine Buchholz and Erwin Marsi. CoNLL-X shared task on multilingual dependency parsing. In *Proceeding of CoNLL-2006*, pages 149–164, New York, NY, 2006.

David Campos, Sérgio Matos, and José Luís Oliveira. Biomedical named entity recognition: a survey of machine-learning tools. *Theory and Applications for Advanced Text Mining*, pages 175–195, 2012.

Xavier Carreras. Experiments with a higher-order projective dependency parser. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 957–961, 2007.

Eugene Charniak. A maximum-entropy-inspired parser. In *Proceedings of the 1st North American chapter of the Association for Computational Linguistics conference*, pages 132–139. Association for Computational Linguistics, 2000.

Danqi Chen and Christopher Manning. A fast and accurate dependency parser using neural networks. In *Proceedings of EMNLP-2014*, pages 740–750, Doha, Qatar, October 2014.

Hao Cheng, Hao Fang, Xiaodong He, Jianfeng Gao, and Li Deng. Bi-directional attention with agreement for dependency parsing. In *Proceedings of EMNLP-2016*, pages 2204–2214, Austin, Texas, November 2016.

Hai Leong Chieu and Hwee Tou Ng. Named entity recognition: a maximum entropy approach using global information. In *Proceedings of CoNLL-2003*, pages 1–7, 2002.

Jason Chiu and Eric Nichols. Named entity recognition with bidirectional lstm-cnns. *Transactions*

*of the Association for Computational Linguistics*, 4:357–370, 2016.

Kevin Clark, Minh-Thang Luong, Christopher D Manning, and Quoc Le. Semi-supervised sequence modeling with cross-view training. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1914–1925, 2018.

Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.

Michael Collins. 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, 2002.

Michael Collins. Head-driven statistical models for natural language parsing. *Computational linguistics*, 29(4):589–637, 2003.

Michael Collins and Yoram Singer. Unsupervised models for named entity classification. In *1999 Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, 1999.

Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537, 2011.

Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What you can cram into a single vector: Probing sentence embeddings for linguistic properties. In *Proceedings ACL-2018 (Volume 1: Long Papers)*, pages 2126–2136, 2018.

Koby Crammer and Yoram Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of machine learning research*, 2(Dec):265–292, 2001.

Koby Crammer and Yoram Singer. Ultraconservative online algorithms for multiclass problems. *Journal of Machine Learning Research*, 3(Jan):951–991, 2003.

Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7(Mar):551–585, 2006.

Leyang Cui and Yue Zhang. Hierarchically-refined label attention network for sequence labeling. In *Proceedings of EMNLP-2019*, pages 4106–4119, 2019.

Hong-Jie Dai, Po-Ting Lai, Yung-Chun Chang, and Richard Tzong-Han Tsai. Enhancing of chemical compound and drug name recognition using representative tag scheme and fine-grained tokenization. *Journal of cheminformatics*, 7(S1):1–10, 2015.

Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, and Stephan Vogel. Understanding and improving morphological learning in the neural machine translation decoder. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 142–151, 2017.

Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. What is one grain of sand in the desert? analyzing individual neurons in deep nlp models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6309–6317, 2019.

Dipanjan Das and Slav Petrov. Unsupervised part-of-speech tagging with bilingual graph-based projections. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 600–609. Association for Computational Linguistics, 2011.

Hal Daumé, John Langford, and Daniel Marcu. Search-based structured prediction. *Machine learning*, 75(3):297–325, 2009.

Yann N Dauphin, Harm de Vries, Junyoung Chung, and Yoshua Bengio. Rmsprop and equilibrated adaptive learning rates for non-convex optimization. *arXiv preprint arXiv:1502.04390*, 2015.

Marie-Catherine De Marneffe, Bill MacCartney, Christopher D. Manning, et al. Generating typed dependency parses from phrase structure parses. In *Proceedings of LREC-2006*, pages 449–454, 2006.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-2019, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.

Cıcero dos Santos, Victor Guimaraes, RJ Niterói, and Rio de Janeiro. Boosting named entity

recognition with neural character embeddings. In *Proceedings of NEWS 2015 The Fifth Named Entities Workshop*, page 25, 2015.

Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.

Finale Doshi-Velez, Mason Kortz, Ryan Budish, Chris Bavitz, Sam Gershman, David O'Brien, Stuart Schieber, James Waldo, David Weinberger, and Alexandra Wood. Accountability of ai under the law: The role of explanation. *arXiv preprint arXiv:1711.01134*, 2017.

Timothy Dozat and Christopher D. Manning. Deep biaffine attention for neural dependency parsing. In *Proceedings of ICLR-2017 (Volume 1: Long Papers)*, Toulon, France, August 2017.

Greg Durrett and Dan Klein. Easy victories and uphill battles in coreference resolution. In *Proceedings of EMNLP-2013*, pages 1971–1982, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.

Chris Dyer, Miguel Ballesteros, Wang Ling, Austin Matthews, and Noah A. Smith. Transition-based dependency parsing with stack long short-term memory. In *Proceedings of ACL-2015 (Volume 1: Long Papers)*, pages 334–343, Beijing, China, July 2015.

Jason M Eisner. Three new probabilistic models for dependency parsing: An exploration. In *Proceedings of COLING-1996 (Volume 1)*, pages 340–345. Association for Computational Linguistics, 1996.

Allyson Ettinger, Ahmed Elgohary, and Philip Resnik. Probing for semantic evidence of composition by means of simple classification tasks. In *Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP*, pages 134–139, 2016.

Oren Etzioni, Michael Cafarella, Doug Downey, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S Weld, and Alexander Yates. Unsupervised named-entity extraction from the web: An experimental study. *Artificial intelligence*, 165(1):91–134, 2005.

Nicolas R Fauceglia, Yiu-Chang Lin, Xuezhe Ma, and Eduard Hovy. Word sense disambiguation via propstore and ontonotes for event mention detection. In *Proceedings of the The 3rd Workshop on EVENTS: Definition, Detection, Coreference, and Representation*, pages 11–15, Denver, Colorado, June 2015.

Radu Florian, Abe Ittycheriah, Hongyan Jing, and Tong Zhang. Named entity recognition through classifier combination. In *Proceedings of HLT-NAACL-2003*, pages 168–171, 2003.

G David Forney. The viterbi algorithm. *Proceedings of the IEEE*, 61(3):268–278, 1973.

Yoav Freund and Robert E Schapire. Large margin classification using the perceptron algorithm. *Machine learning*, 37(3):277–296, 1999.

Yarin Gal and Zoubin Ghahramani. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in Neural Information Processing Systems*, 2016.

J Ganesh, Manish Gupta, and Vasudeva Varma. Interpretation of semantic tweet representations. In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*, pages 95–102, 2017.

Nikhil Garg and James Henderson. Temporal restricted boltzmann machines for dependency parsing. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, pages 11–17. Association for Computational Linguistics, 2011.

Felix A Gers, Nicol N Schraudolph, and Jürgen Schmidhuber. Learning precise timing with lstm recurrent networks. *The Journal of Machine Learning Research*, 3:115–143, 2003.

Rich Caruana Steve Lawrence Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in Neural Information Processing Systems 13: Proceedings of the 2000 Conference*, volume 13, page 402. MIT Press, 2001.

Jesús Giménez and Lluís Màrquez. Svmtool: A general pos tagger generator based on support vector machines. In *In Proceedings of LREC-2004*, 2004.

Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, pages 249–256, 2010.

Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. 2016.

A Graves. Supervised sequence labelling with recurrent neural networks [ph. d. dissertation]. *Technical University of Munich, Germany*, 2008.

Alan Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *Proceedings of ICASSP-2013*, pages 6645–6649. IEEE, 2013.

Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.

Alex Graves and Jürgen Schmidhuber. Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5):602–610, 2005.

Aria Haghighi and Dan Klein. Prototype-driven learning for sequence models. In *Proceedings of the main conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, pages 320–327. Association for Computational Linguistics, 2006.

Jan Hajič, Massimiliano Ciaramita, Richard Johansson, Daisuke Kawahara, Maria Antònia Martí, Lluís Màrquez, Adam Meyers, Joakim Nivre, Sebastian Padó, Jan Štěpánek, et al. The conll-2009 shared task: Syntactic and semantic dependencies in multiple languages. In *Proceedings of CoNLL-2009: Shared Task*, pages 1–18, 2009.

Daniel Hanisch, Katrin Fundel, Heinz-Theodor Mevissen, Ralf Zimmer, and Juliane Fluck. Prominer: rule-based protein and gene entity recognition. *BMC bioinformatics*, 6(1):S14, 2005.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1026–1034, 2015.

James Henderson. Discriminative training of a neural network statistical parser. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 95. Association for Computational Linguistics, 2004.

John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. In *Proceedings of EMNLP-2019*, pages 2733–2743, 2019.

John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of NAACL-2019, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.

Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.

Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenau, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. Robust disambiguation of named entities in text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 782–792. Association for Computational Linguistics, 2011.

Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard H. Hovy, and Eric P. Xing. Harnessing deep neural networks with logic rules. In *Proceedings of ACL-2016*, Berlin, Germany, August 2016.

Zhiheng Huang, Wei Xu, and Kai Yu. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*, 2015.

Richard A Hudson. *Word grammar*. Blackwell Oxford, 1984.

Kevin Humphreys, Robert Gaizauskas, Saliha Azzam, Christian Huyck, Brian Mitchell, Hamish Cunningham, and Yorick Wilks. University of sheffield: Description of the lasie-ii system as used for muc-7. In *Seventh Message Understanding Conference (MUC-7): Proceedings of a Conference Held in Fairfax, Virginia, April 29-May 1, 1998*, 1998.

Mark Johnson. Why doesn't em find good hmm pos-taggers? In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 296–305, 2007.

Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent network architectures. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 2342–2350, 2015.

Daniel Jurafsky and James H Martin. Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition, 2000.

Jerrold J Katz and Jerry A Fodor. The structure of a semantic theory. *language*, 39(2):170–210, 1963.

Ji-Hwan Kim and Philip C Woodland. A rule-based named entity recognition system for speech input. In *Sixth International Conference on Spoken Language Processing*, 2000.

Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Eliyahu Kiperwasser and Yoav Goldberg. Simple and accurate dependency parsing using bidirectional lstm feature representations. *Transactions of the Association for Computational Linguistics*, 4:313–327, 2016.

Arne Köhn. What's in an embedding? analyzing word embeddings through multilingual evaluation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2067–2073, 2015.

Terry Koo and Michael Collins. Efficient third-order dependency parsers. In *Proceedings of ACL-2010*, pages 1–11, Uppsala, Sweden, July 2010.

Terry Koo, Amir Globerson, Xavier Carreras, and Michael Collins. Structured prediction models via the matrix-tree theorem. In *Proceedings of EMNLP-2007*, pages 141–150, Prague, Czech Republic, June 2007.

Terry Koo, Alexander M. Rush, Michael Collins, Tommi Jaakkola, and David Sontag. Dual decomposition for parsing with non-projective head automata. In *Proceedings of EMNLP-2010*, pages 1288–1298, Cambridge, MA, October 2010.

Vijay Krishnan and Christopher D Manning. An effective two-stage model for exploiting non-local dependencies in named entity recognition. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 1121–1128. Association for Computational Linguistics, 2006.

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

GR Krupka and KH IsoQuest. Description of the nerowl extractor system as used for muc-7. In *Proceedings of the 7th Message Understanding Conference, Virginia*, pages 21–28, 2005.

Adhiguna Kuncoro, Miguel Ballesteros, Lingpeng Kong, Chris Dyer, and Noah A. Smith. Distilling an ensemble of greedy dependency parsers into one mst parser. In *Proceedings of EMNLP-2016*, pages 1744–1753, Austin, Texas, November 2016.

Matthieu Labeau, Kevin Löser, Alexandre Allauzen, and Rue John von Neumann. Non-lexical neural architecture for fine-grained pos tagging. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 232–237, 2015.

John Lafferty, Andrew McCallum, and Fernando CN Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of ICML-2001*, volume 951, pages 282–289, 2001.

Yair Lakretz, Germán Kruszewski, Théo Desbordes, Dieuwke Hupkes, Stanislas Dehaene, and Marco Baroni. The emergence of number and syntax units in lstm language models. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 11–20, 2019.

Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. Neural architectures for named entity recognition. In *Proceedings of NAACL-2016*, San Diego, California, USA, June 2016.

Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.

Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.

Tao Lei, Yu Xin, Yuan Zhang, Regina Barzilay, and Tommi Jaakkola. Low-rank tensors for scoring dependency structures. In *Proceedings of ACL-2014 (Volume 1: Long Papers)*, pages 1381–1391, Baltimore, Maryland, June 2014.

Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.

Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering*, 2020.

Shen Li, Joao V Graça, and Ben Taskar. Wiki-ly supervised part-of-speech tagging. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and*

*Computational Natural Language Learning*, pages 1389–1398. Association for Computational Linguistics, 2012.

Wenhui Liao and Sriharsha Veeramachaneni. A simple semi-supervised algorithm for named entity recognition. In *Proceedings of the NAACL HLT 2009 Workshop on Semi-Supervised Learning for Natural Language Processing*, pages 58–65, 2009.

Dekang Lin and Xiaoyun Wu. Phrase clustering for discriminative learning. In *Proceedings of ACL-2009*, pages 1030–1038, 2009.

Wang Ling, Chris Dyer, Alan W Black, and Isabel Trancoso. Two/too simple adaptations of word2vec for syntax problems. In *Proceedings of NAACL-2015*, pages 1299–1304, Denver, Colorado, May–June 2015.

Zachary C Lipton. The mythos of model interpretability. *Queue*, 16(3):31–57, 2018.

Nelson F Liu, Matt Gardner, Yonatan Belinkov, Matthew E Peters, and Noah A Smith. Linguistic knowledge and transferability of contextual representations. In *Proceedings of NAACL-2019, Volume 1 (Long and Short Papers)*, pages 1073–1094, 2019a.

Shifeng Liu, Yifang Sun, Bing Li, Wei Wang, and Xiang Zhao. Hamner: Headword amplified multi-span distantly supervised method for domain specific named entity recognition. *arXiv preprint arXiv:1912.01731*, 2019b.

Xiaohua Liu, Shaodian Zhang, Furu Wei, and Ming Zhou. Recognizing named entities in tweets. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 359–367, 2011.

Yijin Liu, Fandong Meng, Jinchao Zhang, Jinan Xu, Yufeng Chen, and Jie Zhou. Gcdt: A global context enhanced deep transition architecture for sequence labeling. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2431–2441, 2019c.

Yin Lou, Rich Caruana, and Johannes Gehrke. Intelligible models for classification and regression. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 150–158, 2012.

Gang Luo, Xiaojiang Huang, Chin-Yew Lin, and Zaiqing Nie. Joint entity recognition and

disambiguation. In *Proceedings of EMNLP-2015*, pages 879–888, Lisbon, Portugal, September 2015.

Ying Luo, Fengshun Xiao, and Hai Zhao. Hierarchical contextualized representation for named entity recognition. *arXiv preprint arXiv:1911.02257*, 2019.

Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In *Proceedings of EMNLP-2015*, pages 1412–1421, Lisbon, Portugal, September 2015.

Xuezhe Ma and Eduard Hovy. Efficient inner-to-outer greedy algorithm for higher-order labeled dependency parsing. In *Proceedings of the EMNLP-2015*, pages 1322–1328, Lisbon, Portugal, September 2015.

Xuezhe Ma and Eduard Hovy. End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF. In *Proceedings of ACL-2016*, pages 1064–1074, Berlin, Germany, August 2016.

Xuezhe Ma and Eduard Hovy. Neural probabilistic model for non-projective mst parsing. In *Proceedings of IJCNLP-2017 (Volume 1: Long Papers)*, pages 59–69, Taipei, Taiwan, November 2017.

Xuezhe Ma and Fei Xia. Unsupervised dependency parsing with transferring distribution via parallel guidance and entropy regularization. In *Proceedings of ACL-2014*, pages 1337–1348, Baltimore, Maryland, June 2014.

Xuezhe Ma and Hai Zhao. Fourth-order dependency parsing. In *Proceedings of COLING 2012: Posters*, pages 785–796, Mumbai, India, December 2012a.

Xuezhe Ma and Hai Zhao. Probabilistic models for high-order projective dependency parsing. *Technical Report, arXiv:1502.04174*, 2012b.

Xuezhe Ma, Zhengzhong Liu, and Eduard Hovy. Unsupervised ranking model for entity coreference resolution. In *Proceedings of NAACL-2016*, San Diego, California, USA, June 2016.

Xuezhe Ma, Yingkai Gao, Zhiting Hu, Yaoliang Yu, Yuntian Deng, and Eduard Hovy. Dropout with expectation-linear regularization. In *Proceedings of the 5th International Conference on Learning Representations (ICLR-2017)*, Toulon, France, April 2017.

Xuezhe Ma, Zecong Hu, Jingzhou Liu, Nanyun Peng, Graham Neubig, and Eduard Hovy. Stack-pointer networks for dependency parsing. In *Proceedings of ACL-2018 (Volume 1: Long Papers)*, pages 1403–1414, 2018.

Xuezhe Ma, Chunting Zhou, Xian Li, Graham Neubig, and Eduard Hovy. Flowseq: Non-autoregressive conditional sequence generation with generative flow. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4273–4283, 2019.

Christopher D Manning. Part-of-speech tagging from 97% to 100%: is it time for some linguistics? In *Computational Linguistics and Intelligent Text Processing*, pages 171–189. Springer, 2011.

Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993.

Andre Martins, Noah Smith, Mario Figueiredo, and Pedro Aguiar. Dual decomposition with many overlapping components. In *Proceedings of EMNLP-2011*, pages 238–249, Edinburgh, Scotland, UK., July 2011.

Andre Martins, Miguel Almeida, and Noah A. Smith. Turning on the turbo: Fast third-order non-projective turbo parsers. In *Proceedings of ACL-2013 (Volume 2: Short Papers)*, pages 617–622, Sofia, Bulgaria, August 2013.

Hiroshi Maruyama. Structural disambiguation with constraint propagation. In *Proceedings of the 28th annual meeting on Association for Computational Linguistics*, pages 31–38. Association for Computational Linguistics, 1990.

Andrew McCallum and Wei Li. Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 188–191. Association for Computational Linguistics, 2003.

Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. Learned in translation: Contextualized word vectors. In *Advances in Neural Information Processing Systems*, pages 6294–6305, 2017.

Ryan McDonald. *DISCRIMINATIVE LEARNING AND SPANNING TREE ALGORITHMS FOR*

*DEPENDENCY PARSING.* PhD thesis, University of Pennsylvania, 2006.

Ryan McDonald and Joakim Nivre. Analyzing and integrating dependency parsers. *Computational Linguistics*, 37(1):197–230, 2011.

Ryan McDonald and Fernando Pereira. Online learning of approximate dependency parsing algorithms. In *Proceeding of EACL-2006*, 2006.

Ryan McDonald, Koby Crammer, and Fernando Pereira. Online large-margin training of dependency parsers. In *Proceedings of ACL-2005*, pages 91–98, Ann Arbor, Michigan, USA, June 25-30 2005a.

Ryan McDonald, Fernando Pereira, Kiril Ribarov, and Jan Hajic. Non-projective dependency parsing using spanning tree algorithms. In *Proceedings of HLT/EMNLP-2005*, pages 523–530, Vancouver, Canada, October 2005b.

Ryan McDonald, Joakim Nivre, Yvonne Quirmbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Täckström, Claudia Bedini, Núria Bertomeu Castelló, and Jungmee Lee. Universal dependency annotation for multilingual parsing. In *Proceedings of ACL-2013*, pages 92–97, Sofia, Bulgaria, August 2013.

Igor Aleksandrovic Melćuk et al. *Dependency syntax: theory and practice*. SUNY press, 1988.

Andrei Mikheev. A knowledge-free method for capitalized word disambiguation. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 159–166. Association for Computational Linguistics, 1999.

Andrei Mikheev, Marc Moens, and Claire Grover. Named entity recognition without gazetteers. In *Proceedings of the ninth conference on European chapter of the Association for Computational Linguistics*, pages 1–8. Association for Computational Linguistics, 1999.

Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černockỳ, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Eleventh annual conference of the international speech communication association*, 2010.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information*

*processing systems*, pages 3111–3119, 2013.

Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38, 2019.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification. *Lingvisticae Investigationes*, 30(1):3–26, 2007.

David Nadeau, Peter D Turney, and Stan Matwin. Unsupervised named-entity recognition: Generating gazetteers and resolving ambiguity. In *Conference of the Canadian society for computational studies of intelligence*, pages 266–277. Springer, 2006.

Vincent Ng. Supervised noun phrase coreference research: The first fifteen years. In *Proceedings of ACL-2010*, pages 1396–1411, Uppsala, Sweden, July 2010. Association for Computational Linguistics.

Truc-Vien T. Nguyen, Alessandro Moschitti, and Giuseppe Riccardi. Convolution kernels on constituent, dependency and sequential structures for relation extraction. In *Proceedings of EMNLP-2009*, pages 1378–1387, Singapore, August 2009.

J Nivre. An efficient algorithm for projective dependency parsing. 2003. In *Proceedings of the 8th International Workshop on Parsing Technologies (IWPT)*, 2003.

Joakim Nivre. Incremental non-projective dependency parsing. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 396–403, 2007.

Joakim Nivre. Non-projective dependency parsing in expected linear time. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*, pages 351–359. Association for Computational Linguistics, 2009.

Joakim Nivre and Jens Nilsson. Pseudo-projective dependency parsing. In *Proceedings of the*

*43rd Annual Meeting on Association for Computational Linguistics*, pages 99–106. Association for Computational Linguistics, 2005.

Joakim Nivre and Mario Scholz. Deterministic dependency parsing of English text. In *Proceedings of COLING-2004*, pages 64–70, Geneva, Switzerland, August 23-27 2004.

Joakim Nivre, Johan Hall, and Jens Nilsson. Memory-based dependency parsing. In *Proceedings of the Eighth Conference on Computational Natural Language Learning (CoNLL-2004) at HLT-NAACL 2004*, pages 49–56, 2004.

Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. The CoNLL 2007 shared task on dependency parsing. In *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL 2007*, pages 915–932, Prague, Czech Republic, June 2007.

Joakim Nivre, Mitchell Abrams, Željko Agić, Lars Ahrenberg, Lene Antonsen, Maria Jesus Aranzabe, Gashaw Arutie, Masayuki Asahara, Luma Ateyah, Mohammed Attia, et al. Universal dependencies 2.2. 2018.

Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. *arXiv preprint arXiv:1211.5063*, 2012.

Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *Proceedings of ICML-2013*, pages 1310–1318, 2013.

Alexandre Passos, Vineet Kumar, and Andrew McCallum. Lexicon infused phrase embeddings for named entity resolution. In *Proceedings of CoNLL-2014*, pages 78–86, Ann Arbor, Michigan, June 2014.

Nanyun Peng and Mark Dredze. Named entity recognition for chinese social media with jointly trained embeddings. In *Proceedings of EMNLP-2015*, pages 548–554, Lisbon, Portugal, September 2015.

Nanyun Peng and Mark Dredze. Improving named entity recognition for chinese social media with word segmentation representation learning. In *Proceedings of ACL-2016*, Berlin, Germany, August 2016.

Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen-tau Yih. Cross-sentence n-ary relation extraction with graph lstms. *Transactions of the Association for Computational Linguistics*, 5:101–115, 2017.

Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of EMNLP-2014*, pages 1532–1543, Doha, Qatar, October 2014.

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of NAACL-2018, Volume 1 (Long Papers)*, pages 2227–2237, 2018.

Matthew E Peters, Waleed Ammar, Chandra Bhagavatula, and Russell Power. Semi-supervised sequence tagging with bidirectional language models. *arXiv preprint arXiv:1705.00108*, 2017.

Slav Petrov, Dipanjan Das, and Ryan McDonald. A universal part-of-speech tagset. In *Proceedings of LREC-2012*, pages 2089–2096, Istanbul, Turkey, May 2012.

Emily Pitler and Ryan McDonald. A linear-time transition system for crossing interval trees. In *Proceedings of NAACL-2015*, pages 662–671, Denver, Colorado, May–June 2015.

Peng Qian, Xipeng Qiu, and Xuan-Jing Huang. Analyzing linguistic knowledge in sequential model of sentence. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 826–835, 2016a.

Peng Qian, Xipeng Qiu, and Xuan-Jing Huang. Investigating language universal and specific properties in word embeddings. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1478–1488, 2016b.

Alexandra Pomares Quimbaya, Alejandro Sierra Múnera, Rafael Andrés González Rivera, Julián Camilo Daza Rodríguez, Oscar Mauricio Muñoz Velandia, Angel Alberto Garcia Peña, and Cyril Labbé. Named entity recognition over electronic health records through a combined dictionary-based approach. *Procedia Computer Science*, 100(100):55–61, 2016.

Lawrence Rabiner and B Juang. An introduction to hidden markov models. *ieee assp magazine*, 3 (1):4–16, 1986.

Lawrence R Rabiner. A tutorial on hidden markov models and selected applications in speech

121

recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.

Lev Ratinov and Dan Roth. Design challenges and misconceptions in named entity recognition. In *Proceedings of CoNLL-2009*, pages 147–155, 2009.

Yael Ravin and Nina Wacholder. *Extracting names from natural-language text*. Citeseer, 1997.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

Alan Ritter, Sam Clark, Oren Etzioni, et al. Named entity recognition in tweets: an experimental study. In *Proceedings of the conference on empirical methods in natural language processing*, pages 1524–1534. Association for Computational Linguistics, 2011.

Tim Rocktäschel, Michael Weidlich, and Ulf Leser. Chemspot: a hybrid system for chemical named entity recognition. *Bioinformatics*, 28(12):1633–1640, 2012.

Cicero D Santos and Bianca Zadrozny. Learning character-level representations for part-of-speech tagging. In *Proceedings of ICML-2014*, pages 1818–1826, 2014.

Naomi Saphra and Adam Lopez. Understanding learning dynamics of language models with svcca. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3257–3267, 2019.

Satoshi Sekine and Chikashi Nobata. Definition, dictionaries and tagger for extended named entity hierarchy. In *LREC*. Lisbon, Portugal, 2004.

Burr Settles. Biomedical named entity recognition using conditional random fields and rich feature sets. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications (NLPBA/BioNLP)*, pages 107–110, 2004.

Petr Sgall, Eva Hajicová, Eva Hajicová, Jarmila Panevová, and Jarmila Panevova. *The meaning of the sentence in its semantic and pragmatic aspects*. Springer Science & Business Media, 1986.

Rahul Sharnagat. Named entity recognition: A literature survey. *Center For Indian Language Technology*, 2014.

Libin Shen, Giorgio Satta, and Aravind Joshi. Guided learning for bidirectional sequence classification. In *Proceedings of ACL-2007*, volume 7, pages 760–767, 2007.

Xing Shi, Inkit Padhi, and Kevin Knight. Does string-based neural mt learn source syntax? In *Proceedings of EMNLP-2016*, pages 1526–1534, 2016.

Aaron Smith, Miryam de Lhoneux, Sara Stymne, and Joakim Nivre. An investigation of the interactions between pre-trained word embeddings, character models and POS tags in dependency parsing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2711–2720, Brussels, Belgium, October-November 2018.

David A. Smith and Noah A. Smith. Probabilistic models of nonprojective dependency trees. In *Proceedings of EMNLP-2007*, pages 132–140, Prague, Czech Republic, June 2007.

Anders Søgaard. Semi-supervised condensed nearest neighbor for part-of-speech tagging. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 48–52, Portland, Oregon, USA, June 2011.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014a.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014b.

Karl Stratos, Michael Collins, and Daniel Hsu. Unsupervised part-of-speech tagging with anchor hidden markov models. *Transactions of the Association for Computational Linguistics*, 4: 245–257, 2016.

Xu Sun. Structure regularization for structured prediction. In *Advances in Neural Information Processing Systems*, pages 2402–2410, 2014.

Mihai Surdeanu, Richard Johansson, Adam Meyers, Lluís Màrquez, and Joakim Nivre. The conll-2008 shared task on joint parsing of syntactic and semantic dependencies. In *Proceedings of CoNLL-2008*, pages 159–177, 2008.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.

Oscar Täckström, Dipanjan Das, Slav Petrov, Ryan McDonald, and Joakim Nivre. Token and type constraints for cross-lingual part-of-speech tagging. *Transactions of the Association for Computational Linguistics*, 1:1–12, 2013.

Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks. In *Proceedings ACL-2015 (Volume 1: Long Papers)*, pages 1556–1566, Beijing, China, July 2015.

Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. olmpics–on what language model pre-training captures. *arXiv preprint arXiv:1912.13283*, 2019.

Ben Taskar, Carlos Guestrin, and Daphne Koller. Max-margin markov networks. In *Advances in neural information processing systems*, pages 25–32, 2004.

Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel Bowman, Dipanjan Das, et al. What do you learn from context? probing for sentence structure in contextualized word representations. In *Proceedings of the 7th International Conference on Learning Representations (ICLR-2019)*, 2019.

Ivan Titov and James Henderson. A latent variable model for generative dependency parsing. In *Proceedings of the Tenth International Conference on Parsing Technologies*, pages 144–155, Prague, Czech Republic, June 2007a.

Ivan Titov and James Henderson. Fast and robust multilingual dependency parsing with a generative latent variable model. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 947–951, 2007b.

Erik F. Tjong Kim Sang and Fien De Meulder. Introduction to the conll-2003 shared task: Language-independent named entity recognition. In *Proceedings of CoNLL-2003 - Volume 4*, pages 142–147, Stroudsburg, PA, USA, 2003.

Erik F. Tjong Kim Sang and Jorn Veenstra. Representing text chunks. In *Proceedings of EACL'99*, pages 173–179. Bergen, Norway, 1999.

Kentaro Torisawa et al. Exploiting wikipedia as external knowledge for named entity recognition. In *Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL)*, pages 698–707, 2007.

Kristina Toutanova, Dan Klein, Christopher D Manning, and Yoram Singer. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of NAACL-HLT-2003, Volume 1*, pages 173–180, 2003.

Ke M Tran, Arianna Bisazza, and Christof Monz. The importance of being recurrent for modeling hierarchical structure. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4731–4736, 2018.

William Thomas Tutte. *Graph theory*, volume 11. Addison-Wesley Menlo Park, 1984.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in Neural Information Processing Systems*, pages 2692–2700, 2015.

Elena Voita and Ivan Titov. Information-theoretic probing with minimum description length. *arXiv preprint arXiv:2003.12298*, 2020.

Elena Voita, Rico Sennrich, and Ivan Titov. The bottom-up evolution of representations in the transformer: A study with machine translation and language modeling objectives. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4387–4397, 2019.

Ekaterina Vylomova, Trevor Cohn, Xuanli He, and Gholamreza Haffari. Word representation models for morphologically rich languages in neural machine translation. In *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, pages 103–108, 2017.

Wenhui Wang and Baobao Chang. Graph-based dependency parsing with bidirectional lstm. In *Proceedings of ACL-2016 (Volume 1: Long Papers)*, pages 2306–2315, Berlin, Germany,

August 2016.

Yu-Hsuan Wang, Cheng-Tao Chung, and Hung-yi Lee. Gate activation signal analysis for gated recurrent neural networks and its correlation with phoneme boundaries. *arXiv preprint arXiv:1703.07588*, 2017.

David Weiss, Chris Alberti, Michael Collins, and Slav Petrov. Structured training for neural network transition-based parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 323–333, 2015a.

David Weiss, Chris Alberti, Michael Collins, and Slav Petrov. Structured training for neural network transition-based parsing. In *Proceedings of ACL-2015 (Volume 1: Long Papers)*, pages 323–333, Beijing, China, July 2015b.

Yorick Wilks. A preferential, pattern-seeking, semantics for natural language inference. *Artificial intelligence*, 6(1):53–74, 1975.

Zhizheng Wu and Simon King. Investigating gated recurrent networks for speech synthesis. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5140–5144. IEEE, 2016.

Jun Xie, Haitao Mi, and Qun Liu. A novel dependency-to-string model for statistical machine translation. In *Proceedings of EMNLP-2011*, pages 216–226, Edinburgh, Scotland, UK., July 2011.

Nianwen Xue, Fu-Dong Chiou, and Martha Palmer. Building a large-scale annotated chinese corpus. In *Proceedings of COLING-2002*, pages 1–8, 2002.

Hiroyasu Yamada and Yuji Matsumoto. Statistical dependency analysis with support vector machines. In *Proceedings of IWPT*, volume 3, pages 195–206. Nancy, France, 2003.

Dani Yogatama, Cyprien de Masson d'Autume, Jerome Connor, Tomas Kocisky, Mike Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, et al. Learning and evaluating general linguistic intelligence. *arXiv preprint arXiv:1901.11373*, 2019.

Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*,

2012.

Hao Zhang and Ryan McDonald. Enforcing structural diversity in cube-pruned dependency parsing. In *Proceedings of ACL-2014 (Volume 2: Short Papers)*, pages 656–661, Baltimore, Maryland, June 2014.

Kelly Zhang and Samuel Bowman. Language modeling teaches you more than translation does: Lessons learned through auxiliary syntactic task analysis. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 359–361, 2018.

Shaodian Zhang and Noémie Elhadad. Unsupervised biomedical named entity recognition: Experiments with clinical and biological texts. *Journal of biomedical informatics*, 46(6): 1088–1098, 2013.

Xingxing Zhang, Jianpeng Cheng, and Mirella Lapata. Dependency parsing as head selection. *arXiv preprint arXiv:1606.01280*, 2016.

Yuan Zhang, Tao Lei, Regina Barzilay, and Tommi Jaakkola. Greed is good if randomized: New inference for dependency parsing. In *Proceedings of EMNLP-2014*, pages 1013–1024, Doha, Qatar, October 2014.

Yue Zhang and Joakim Nivre. Transition-based dependency parsing with rich non-local features. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 188–193, Portland, Oregon, USA, June 2011.

GuoDong Zhou and Jian Su. Named entity recognition using an hmm-based chunk tagger. In *proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pages 473–480. Association for Computational Linguistics, 2002.