

# Neural dependency parsing

(Parsing zależnościowy za pomocą sieci neuronowej)

Michał Zapotoczny

Praca magisterska

**Promotor:** dr Jan Chorowski

Uniwersytet Wrocławski  
Wydział Matematyki i Informatyki  
Instytut Informatyki

31 lipca 2017

Michał Zapotoczny

.....

.....

(adres zameldowania)

.....

.....

(adres korespondencyjny)

PESEL: .....

e-mail: .....

Wydział Matematyki i Informatyki

stacjonarne studia II stopnia

kierunek: informatyka

nr albumu: 248100

### **Oświadczenie o autorskim wykonaniu pracy dyplomowej**

Niniejszym oświadczam, że złożoną do oceny pracę zatytułowaną *Neural dependency parsing* wykonałem samodzielnie pod kierunkiem promotora, dr. Jana Chorońskiego. Oświadczam, że powyższe dane są zgodne ze stanem faktycznym i znane mi są przepisy ustawy z dn. 4 lutego 1994 r. o prawie autorskim i prawach pokrewnych (tekst jednolity: Dz. U. z 2006 r. nr 90, poz. 637, z późniejszymi zmianami) oraz że treść pracy dyplomowej przedstawionej do obrony, zawarta na przekazanym nośniku elektronicznym, jest identyczna z jej wersją drukowaną.

Wrocław, 31 lipca 2017

(czytelny podpis)

## **Abstract**

We present a dependency parser, implemented as single neural network which works directly on character representation of the words. Unlike typical approaches our model doesn't require any preprocessing of the data, like part-of-speech tagging. Additionally we present models trained jointly on two languages, which can obtain state-of-the art results.

---

Prezentujemy parser zależnościowy zaimplementowany jako pojedyncza sieć neuronowa, która działa bezpośrednio na literach słów. W przeciwieństwie do typowych sposobów parsingu nasz model nie wymaga żadnego pre-processingu danych, (w szczególności tagowania słów). Dodatkowo prezentujemy modele wytrenowane na dwóch językach jednocześnie, które uzyskują wyniki konkurencyjne do najlepszych dostępnych parserów.



# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                                     | <b>7</b>  |
| 1.1      | Dependency parsing . . . . .                            | 8         |
| 1.1.1    | Parsing algorithms . . . . .                            | 9         |
| 1.1.2    | Evaluation of parsing algorithms . . . . .              | 10        |
| 1.1.3    | Universal Dependencies . . . . .                        | 10        |
| 1.2      | Neural Networks . . . . .                               | 12        |
| <b>2</b> | <b>Neural dependency parser</b>                         | <b>13</b> |
| 2.1      | Overview of the network architecture . . . . .          | 13        |
| 2.1.1    | Reader . . . . .  | 13        |
| 2.1.2    | Tagger . . . . .  | 15        |
| 2.1.3    | Parser . . . . .  | 15        |
| 2.2      | Training . . . . .                                      | 16        |
| 2.2.1    | Training criterion . . . . .                            | 16        |
| 2.2.2    | Regularization, optimizer and hyperparameters . . . . . | 17        |
| 2.3      | Multilingual training . . . . .                         | 17        |
| <b>3</b> | <b>Results</b>  | <b>19</b> |
| 3.1      | Single language . . . . .                               | 19        |
| 3.1.1    | Impact of <i>reader</i> and <i>predictor</i> . . . . .  | 20        |
| 3.1.2    | Impact of decoding algorithm . . . . .                  | 20        |
| 3.1.3    | Recurrent state size . . . . .                          | 22        |
| 3.1.4    | Word pieces model . . . . .                             | 22        |

|          |  |           |
|----------|--|-----------|
| 3.2      | Multilanguage training . . . . .         | 24        |
| 3.2.1    | Polish-Polish multilang parser . . . . . | 25        |
| <b>4</b> | <b>Summary</b>                           | <b>27</b> |
|          | <b>Bibliography</b>                      | <b>29</b> |

# Chapter 1

## Introduction

The ability to communicate with the user in a natural language is a major driving force in development of natural language processing algorithms. One of the basic tasks in a NLP pipeline is parsing by which we can describe structure of sentences. There exist two basic parsing techniques: constituency and dependency parsing. With a constituency parser we break the sentence into phrases, which can be further broken into smaller sub-phrases. Example of such parsing is shown in Figure 1.1.

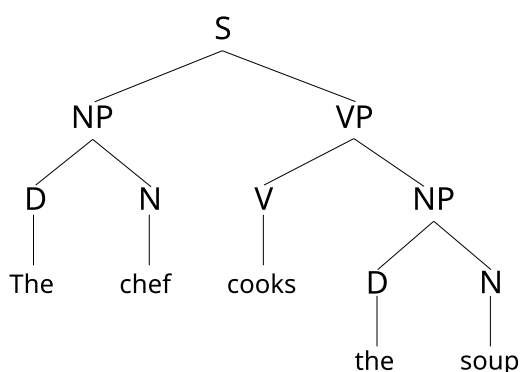


Figure 1.1: A sample constituency parse tree

In dependency parsing each word (called the *dependent*) is connected via a labelled arc to another word of the sentence (called the *head*) or to the special *ROOT* vertex, forming a directed tree. Example of such parsing is depicted in Figure 1.2.

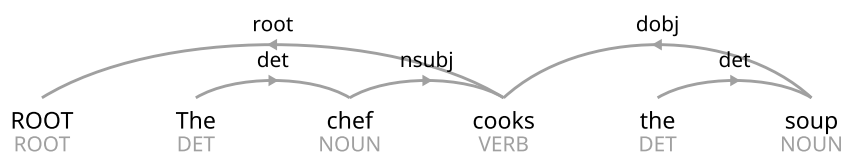


Figure 1.2: A sample dependency parse tree

## 1.1 Dependency parsing

A major advantage of dependency parsers over constituency ones are their independence from the word ordering. It is important in morphologically rich languages like Polish or Czech where the ordering is very flexible, and thus related words can be far apart from each other. Additionally the head-dependent relationship is a good approximation to semantic relationship between words [Covington, 2001] which is important for tasks like question answering or information extraction.

The labels of head-dependent arcs tells us about grammatical function that dependent word have in respect to the head. In table 1.1 we show some of the most popular labels for the English language.

| Label    | Description         | Example                                      |
|----------|---------------------|--|
| CASE     | Case marking        | <i>From</i> Friday 's Daily <b>Star</b>      |
| NSUBJ    | Nominal subject     | <i>Musharraf</i> <b>calls</b> the bluff      |
| NMOD     | Nominal modifier    | India <b>defensive</b> over Sri <i>Lanka</i> |
| DET      | Determiner          | That he missed <i>a</i> <b>physical</b> ?    |
| DOBJ     | Direct object       | Did you <b>know</b> <i>that</i> ?            |
| ADVMOD   | Adverbial modifier  | <i>So</i> Bush <b>stopped</b> flying .       |
| AMOD     | Adjectival modifier | Six weeks of <i>basic</i> <b>training</b> .  |
| COMPOUND | Compound            | Bush 's <i>National</i> <b>Guard</b> years   |

Table 1.1: Most popular labels from English treebank. Dependents are *italic*, while heads are **bold**

### Projectivity

We can distinguish two types of dependency trees: projective and nonprojective. We say that head-dependent arc is projective if there exists a path from head to every word between head and dependent in the sentence. A dependency tree is projective if all its arcs are projective. Intuitively we can say that nonprojective trees are trees that cannot be drawn without crossing the edges (see Figure 1.3). Although English dataset has less than 5% of nonprojective sentences other languages can have significant amount - like Czech (12%) or Ancient Greek (63%) [Straka et al., 2015]. It is important to acknowledge this fact because some of the presented parsing methods can only produce projective trees.

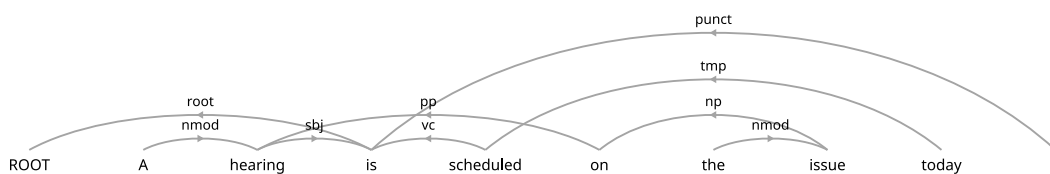


Figure 1.3: A sample nonprojective dependency parse tree. The  $\text{on} \rightarrow \text{hearing}$  arc is nonprojective and thus the whole tree becomes nonprojective.



### 1.1.1 Parsing algorithms

In the following section we will present two basic approaches to dependency parsing. Although they differ in complexity and output flexibility, both of them depend on supervised machine learning techniques.

#### Transition based

A configuration  $c$  is a triple  $(s, b, A)$  containing stack  $s$ , buffer  $b$  and set of dependency arcs  $A$ . In transition based parsing we find a transition sequence from an initial configuration to a terminal one. For given sentence  $w_1, \dots, w_n$  the initial configuration is  $([\text{ROOT}], [w_1, \dots, w_n], \emptyset)$ . and the terminal one is  $(\text{ROOT}, [], A)$  in which  $A$  is a resulting parse tree. There are three possible transitions:

- **LEFTARC(1)** - Pop from the stack elements  $s_1$  and  $s_2$ , add arc  $s_1 \rightarrow s_2$  with label  $l$  to the set  $A$  and push  $s_1$  back to the stack
- **RIGHTARC(1)** - Pop from the stack elements  $s_1$  and  $s_2$ , add arc  $s_2 \rightarrow s_1$  with label  $l$  to the set  $A$  and push  $s_2$  back to the stack
- **SHIFT** - Move first element from the buffer to the stack

An example of parsing sequence is shown in table 1.2.

| Stack                   | Buffer                | Transition     | New arcs            |
|-------------------------|-----------------------|----------------|---------------------|
| [ROOT]                  | [He has good control] | SHIFT          | -                   |
| [ROOT He]               | [has good control]    | SHIFT          | -                   |
| [ROOT He has]           | [good control]        | LEFTARC(nsubj) | nsubj(has, He)      |
| [ROOT has]              | [good control]        | SHIFT          | -                   |
| [ROOT has good]         | [control]             | SHIFT          | -                   |
| [ROOT has good control] | []                    | LEFTARC(amod)  | amod(control, good) |
| [ROOT has control]      | []                    | RIGHTARC(dobj) | dobj(has, control)  |
| [ROOT has]              | []                    | RIGHTARC(root) | root(ROOT, has)     |
| [ROOT]                  | []                    | -              | -                   |

Table 1.2: Transitions needed to parse sentence "He has good control"

Because every word of the sentence has to be shifted to the stack exactly once and ARC operations reduce stack size by one we will have exactly  $2n$  transitions from initial to terminal configuration. Thus the time complexity of this algorithm is  $O(n)$  where  $n$  is number of words in a sentence. Using this algorithm we can represent any projective tree [Nivre, 2008] (nonprojectives cannot be represented).

The decision which transition to use between consecutive configurations is obtained through machine learning techniques. This can vary from simple linear SVM [Nivre et al., 2005] to neural networks [Chen and Manning, 2014, Andor et al., 2016].

## Graph based

Alternative approach is to use graph-based algorithms. The idea is that we define a space of candidate arcs, find a model that scores each arc and then use some tree building algorithm to find a final dependency tree with the highest score.

The candidate space and scoring depends on the used learning model. There exists many parsing algorithms, but one of the most popular is Chu-Liu-Edmonds [Edmonds, 1966]. The Chu-Liu-Edmonds algorithm, given a directed graph produces minimum spanning arborescence (minimum spanning tree for directed graph) - it allows us to produce any parse tree (including nonprojective ones).

Formally we present this method as Algorithm 1, but the general idea is that each vertex chooses an incoming edge of minimum weight. If this graph does not contain any cycle then we are done, else we break the cycle in place that would give us the best result.

### 1.1.2 Evaluation of parsing algorithms

Having predicted dependency tree for a particular sentence we have to be able to compare it with the gold-standard parse. There are two main evaluation metrics:

- **Unlabelled Attachment Score (UAS)**, is a percentage of words with correct predicted head
- **Labelled Attachment Score (LAS)**, is a percentage of words with correct predicted head **and** correct predicted label

In the experimental section we will show both scores, whenever available.

### 1.1.3 Universal Dependencies

Up till recently most development of dependency parsers was done in single language setting, where we have different parser for every language we want to use, each using language-specific morphosyntactical features. The Universal Dependencies project [Nivre et al., 2015] aims to provide a cross-linguistically consistent treebank annotation. It is based on previous work on universal Stanford Dependencies [De Marneffe et al., 2006], Google universal pos tags [Petrov et al., 2011] and the Interset interlingua [Zeman, 2008]. The UD project unifies part of speech tags and dependency relations labels for 30+ languages under common CONNL-U format. The main advantage is that we can use the same parsing system for every language (data has consistent format) additionally because the morphosyntactical data has common format we can try to combine data coming from different languages (see section 2.3).

**Data:** Directed graph  $G = (V, E)$ , root node  $r \in V$ , weight function  $w(e), e \in E$

**Result:** Spanning arborescence of minimum weight rooted at  $r$

**Function** *edmonds*( $G = (V, E), r, w$ )

- Remove any edge from  $G$  whose destination is  $r$  ;
- For each  $v \in V - \{r\}$  find incoming edge of minimum weight, denote it as  $\pi(v)$  ;
- Denote new set of edges as  $P = (\pi(v), v) | v \in V - \{r\}$  ;
- if**  $P$  does not contain a cycle **then**
  - | **return**  $P$
- else**
  - Choose any cycle from  $P$  and denote it as  $C = (C_v, C_e)$  ;
  - Define a new directed graph  $G' = (V', E')$ , and weighting function  $w'$ , where  $V' = V - C_v + v_c$  ;
  - Define  $E'$  and  $w'$  as follows: **for**  $(u, v) \in E$  **do**
    - if**  $u \notin C_e, v \in C_e$  **then**
      - |  $E' = E' + (u, v_c), w'(u, v_c) = w(u, v) - w(\pi(v), v)$
    - else if**  $u \in C_e, v \notin C_e$  **then**
      - |  $E' = E' + (v_c, v), w'(v_c, v) = w(u, v)$
    - else if**  $u \notin C_e, v \notin C_e$  **then**
      - |  $E' = E' + (u, v), w'(u, v) = w(u, v)$
    - end**
  - For each new edge remember its origin ;
  - end**
  - Denote  $A' = \text{edmonds}(G', r, w')$  ;
  - Let  $(u, v_c)$  be (unique) incoming edge to  $v_c$  in  $A'$  ;
  - Let  $(u, v) \in E(v \in C)$  be a corresponding edge from  $G$  ;
  - Remove edge  $(\pi(v), v)$  from  $C$  (breaking the cycle) ;
  - Each remaining edge in  $C$  together with edges from  $E$  that are also present in  $A'$  forms a minimum spanning arborescence. ;
- end**

**Algorithm 1:** Chu-Liu-Edmonds Algorithm. The standard algorithm has  $O(EV)$  running time, but there exist faster implementation by [Gabow et al., 1986] with  $O(E + V \log V)$  running time.

## 1.2 Neural Networks

A neural network is a mathematical model loosely based on how the brain works. The models were proven to be very flexible and obtained state of the art results in many tasks like machine translation [Bahdanau et al., 2014], caption generation [Xu et al., 2015] or image recognition [Szegedy et al., 2014]. A neural network is defined by its topology, activation functions and parameters. The first two are chosen by the author, while the parameters are obtained by a learning procedure.

The topology of the network is telling us how the layers of neurons are connected to each other. Each layer has a number of neurons which have some incoming indices  $i$ , associated activation function  $\phi$  and some parameters  $w$ . The neuron combines the incoming signals, usually by taking weighted average and then applying the activation function:  $\phi(\sum_{k=0}^{|i|} w_k i_k)$ . This value is then fed to the connected neurons from the next layer.

The activation function can have any form, but the most common are hyperbolic tangent  $\tanh$  and ReLU [Nair and Hinton, 2010].

By the learning procedure we understand an optimization process that given a loss function  $L$ , data  $x$  and some parameters  $W$  to tune, tries to minimize  $L(x)$ . Because usually the data  $x$  size is big we cannot use simple gradient descent algorithm on all data at once because it would not fit into the memory. Instead we are using stochastic gradient descent to train only on small subset of data at once (this subset is called a minibatch).

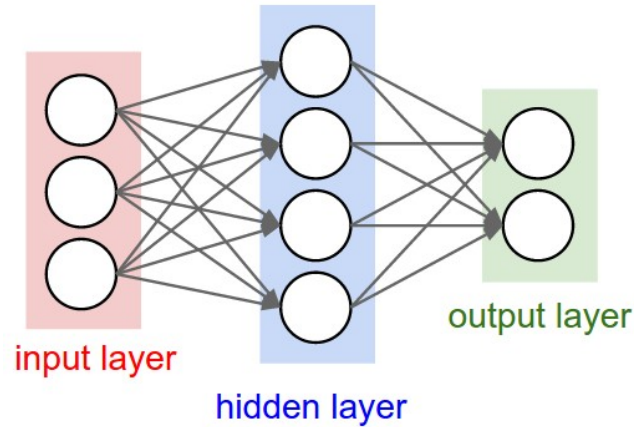


Figure 1.4: An sample neural network architecture

## Chapter 2

# Neural dependency parser

In this chapter we will present the overall topology of our neural network. Additionally we will describe the training procedure and how the proper parameters were selected.

### 2.1 Overview of the network architecture

The network architecture consists of three main parts: reader, tagger and parser (see Figure 2.1). The reader subnetwork is evaluated on each individual word in a sentence, and using convolutions on their orthographic representation produces words embeddings. Next, the tagger subnetwork implemented as bidirectional RNN equips each word with a context of whole sentence. Finally parser part computes dependency tree parent for each word using attention mechanism [Vinyals et al., 2015] after which network computes appropriate dependency label. In the following paragraphs we will describe all parts in detail.

#### 2.1.1 Reader

As stated before the reader subnetwork is run on each word producing its embedding. This architecture is based on [Kim et al., 2015]. Each word  $w$  is represented by sequence of its characters plus a special beginning-of-word **<bow>** and end-of-word **<eow>** tokens. First, we find low-dimensional characters embeddings which are concatenated to form a matrix  $C^w$ . Then we run 1D convolutional filters on  $C^w$  which then is reduced to vector of filter responses computed as:

$$R_i^w = \max(C^w * F^i) \quad (2.1)$$

Where  $F^i$  is  $i$ -th filter. The purpose of the convolutions is to react to specific part of words (because we have **<bow>** and **<eow>** tokens it can also react to prefixes and suffixes) which in morphologically rich languages such as Polish can depict its grammatical role.

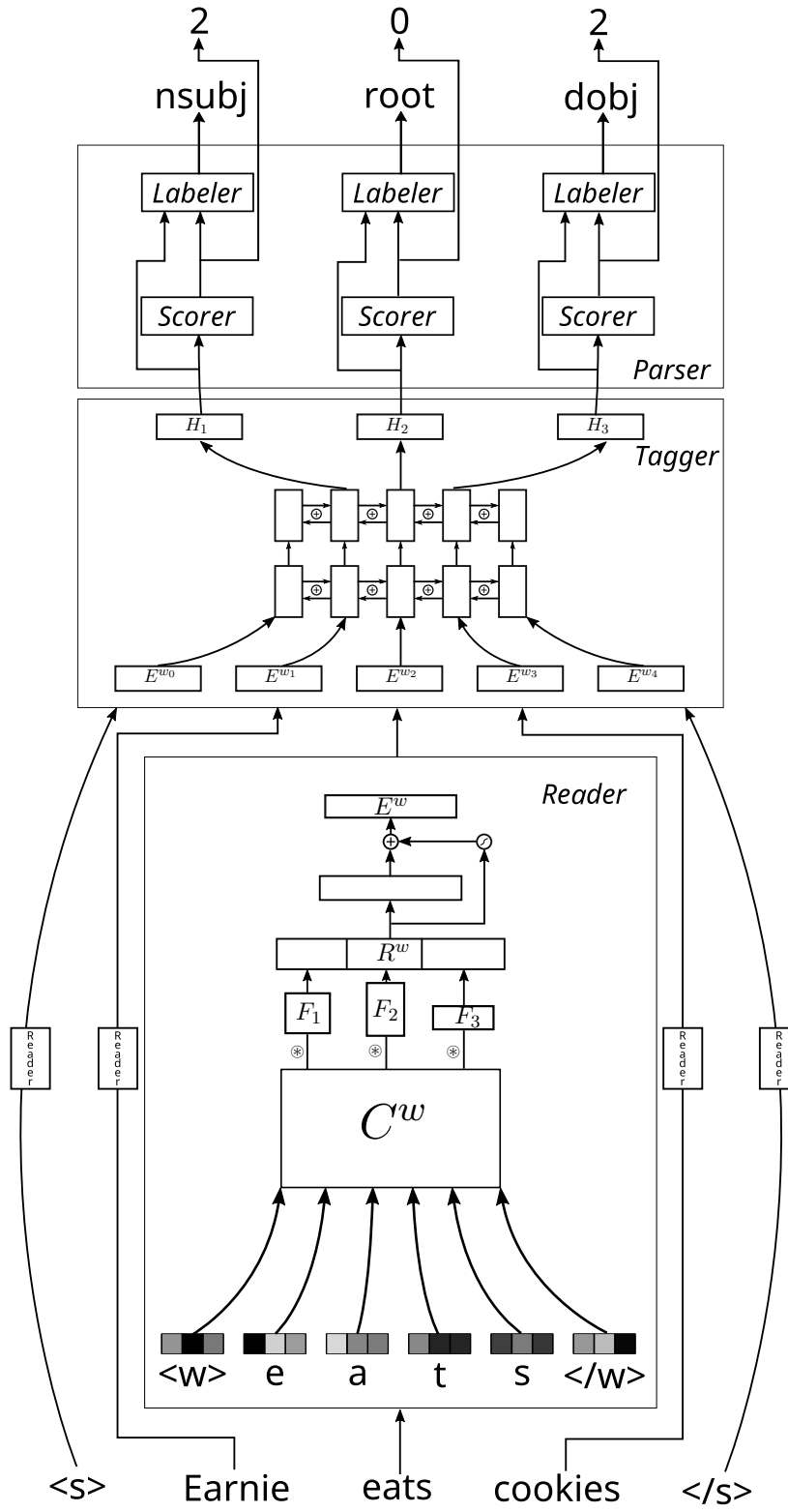


Figure 2.1: The model architecture.

Finally we transform filter responses  $R^w$  with a simple multi-layer perceptron <sup>1</sup> obtaining the final word embedding  $E^w = \text{MLP}(R^w)$ .

### 2.1.2 Tagger

Having obtained the word embeddings  $E^w$  we can proceed with actually "reading" whole sentence. To do this we use multi-layer bidirectional RNN ([Schuster and Paliwal, 1997]) (we have evaluated LSTM [Hochreiter and Schmidhuber, 1997] and GRU [Cho et al., 2014] types). We combine the backward and forward passes by adding them.

#### POS tag predictor

To prevent overfitting and encourage network to compute morphological features we can add additional training objective. It works by taking output from one of the tagger BiRNN layers and use it to predict available part-of-speech tags for each word. The result is not feeded back to the rest of the network because we think it would introduce too much noise.

### 2.1.3 Parser

The last part of the network serves two purposes: to find head for each word and to compute label for that edge.

#### Finding the head word

We use a method similar to [Vinyals et al., 2015]. The input to this part are word annotations  $H_0, H_1, \dots, H_n$  (where  $H_0$  serves as a root word) produced by the tagger. For each of the words  $1, 2, \dots, n$  we compute probability distribution which tell us where the head of current word should be  $(0, 1, 2, \dots, n)$ . This computation (called *scorer*) is implemented as small feedforward network  $s(w, l) = f(H_w, H_l)$ , where  $w \in 1, 2, \dots, n, l \in 0, 1, 2, \dots, n$ .

#### Finding edge label

Output of the scorer can already be interpreted as pointer network, but in order to use it as attention for computing dependency label we have to normalize it:

$$p(w, l) = \frac{f(H_w, H_l)}{\sum_{i=0}^n f(H_w, H_i)} \quad (2.2)$$

---

<sup>1</sup>Which are just linear transformations followed by non-linearity

The dependency label is computed by small Maxout network [Goodfellow et al., 2013], which takes the current word annotation  $H_w$  and heads annotation  $A_w$ . This part is called the *labeller*. We investigated two variations of head annotation  $A_w$

#### *Soft attention*

Which is a weighted average of normalized attention 2.2 and words annotations  $H$ .

$$A_w = \sum_{i=0}^n p(w, i) H_i$$

#### *Hard attention*

Here we only use a head annotation.

$$A_w = H_h$$

During training we use ground-truth head location, whereas during evaluation we use head word computed in previous step.

### Decoding algorithm

The *scorer* give us a  $n \times n + 1$  matrix of head dependency preference for each word. From this matrix we have find a set of dependencies that satisfy some constraints (exactly one root dependant, no cycles). We investigated two possibilities for such decoding: a greedy algorithm, and Chu-Liu-Edmonds [Edmonds, 1966].

## 2.2 Training

### 2.2.1 Training criterion

Every neural network needs a training criterion which will be optimized. Here we have 3 individual training criterion combined together as linear combination. Those are:

- The negative log-likelihood loss  $L_h$  on finding the proper head for each word. The training signal is propagated from scorer down to the reader.
- The negative log-likelihood loss  $L_l$  on finding the dependency label. With a soft attention it is propagated through the whole network (excluding pos tagging part), while with the hard attention we do not propagate through the scorer.
- The (optional) negative log-likelihood loss  $L_t$  on predicting pos tags. This error is backpropagated only through pos-predictor and part of tagger down to the reader.



The final loss is:

$$L = \alpha_h L_h + \alpha_l L_l + \alpha_t L_t \quad (2.3)$$

### 2.2.2 Regularization, optimizer and hyperparameters

Regularization is an essential part of neural network training. It improves generalization and prevents overfitting. One of the most popular regularization technique is called Dropout [Srivastava et al., 2014]. Using it, we randomly drop part of the connections from a certain network layer during training. In our case dropout is applied to the *reader* output, between the BiRNN layers of the *tager* and to the *labeller*.

The models are trained using Adedelta [Zeiler, 2012] learning rule, with weight decay and adaptive gradient clipping [Chorowski et al., 2014]. All experiments are early stopped on validation set UAS.

Hyperparameter selection is crucial for neural networks to obtain good results. For example, comparing our first experiments with architecture depicted above to the best single-language results (using the same basic architecture) gave us about 5% boost in UAS score on Polish language.

To find the best hyperparameters we have used the Spearmint system [Snoek et al., 2012] invoked on polish dataset. The chosen parameters are as follows. The *reader* embeds each character into 15 dimensions, and contains 1050 filters (50-k filters of length k for  $k = 1, 2, \dots, 6$ ) whose outputs are projected into 512 dimensions transformed by a 3 equally sized layers of feedforward neural network with ReLU activation. The *tager* contains 2 BiRNN layers of GRU units with 548 hidden states for both forward and backward passes which are later aggregated using addition. Therefore the hidden states of the *tager* are also 548-dimensional. The *POS tag predictor* consists of a single affine transformation followed by a SoftMax predictor for each POS category. The *scorer* uses a single layer of 384 tanh for head word scoring while the *labeller* uses 256 Maxout units (each using 2 pieces) to classify the relation label [Goodfellow et al., 2013]. The training cost used the constants  $\alpha_h = 0.6, \alpha_l = 0.4, \alpha_t = 1.0$ . We apply 20% dropout to the *reader* output, 70% between the BiRNN layers of the *tager* and 50% to the *labeller*. The weight decay is 0.95.

## 2.3 Multilingual training

The big problem of learning to parse is small number of gold standard dependency trees for many languages (including Polish). With small number of examples neural networks do not generalize well and can more easily overfit. There also exist languages with good, standardized treebank like Czech Prague Treebank [Bejček et al., 2013].

Because our model is purely neural network we can incorporate a multitask learning [Caruana, 1997]. It allows the network to learn multiple tasks at the same time, sharing common patterns and distinguishing differences. Additionally, because we are using Universal Dependencies treebanks [Nivre et al., 2015] we have common standardized format across many languages, which allowed for easy implementation of multitask learning and experiments across different languages.

The multilingual model for  $n$  languages can be viewed as  $n$  copies of our basic model, but sharing part of the parameters. To unify input/output for each of the models we sum possible outputs for each data category (characters, POS categories, dependency labels). If some category doesn't exist within a particular language, we use a special UNK token. To actually make use of multitask learning we must share at least part of the parameters of all models. We experimented with different sharing strategies, from share-everything to only sharing the *parser* part. Additionally to prevent over-representation of some languages during training we sample (on each epoch) only portion of the available data so that each language have equal number of examples (equal to number of samples of the smallest language).

## Chapter 3

# Results

### 3.1 Single language

All experiments depicted in this section are based on a configuration shown in section 2.2.2. In Table 3.1 we show our baseline results for a subset of the UD languages (due to limited computational resources we were not able to train models for all of them). The results are compared to SyntaxNet [Andor et al., 2016] and ParseySaurus [Alberti et al., 2017], both from Google. The ParseySaurus is based on SyntaxNet, and uses characters as input, while the original SyntaxNet uses word embeddings. Both are transition-based, and need no external POS tagger. It is worth noting that SyntaxNet has different hyperparameter configuration for each language, while our parser uses the same configuration across all languages.

| language      | #sentences | Ours         |              | SyntaxNet |       | ParseySaurus |              |
|---------------|------------|--------------|--------------|-----------|-------|--------------|--------------|
|               |            | UAS          | LAS          | UAS       | LAS   | UAS          | LAS          |
| Czech         | 87 913     | <b>91.41</b> | <b>88.18</b> | 89.47     | 85.93 | 89.09        | 84.99        |
| Polish        | 8 227      | 90.26        | 85.32        | 88.30     | 82.71 | <b>91.86</b> | <b>87.49</b> |
| Russian       | 5 030      | 83.29        | 79.22        | 81.75     | 77.71 | <b>84.27</b> | <b>80.65</b> |
| German        | 15 892     | 82.67        | 76.51        | 79.73     | 74.07 | <b>84.12</b> | <b>79.05</b> |
| English       | 16 622     | 87.44        | 83.94        | 84.79     | 80.38 | <b>87.86</b> | <b>84.45</b> |
| French        | 16 448     | <b>87.25</b> | <b>83.50</b> | 84.68     | 81.05 | 86.61        | 83.1         |
| Ancient Greek | 25 251     | <b>78.96</b> | <b>72.36</b> | 68.98     | 62.07 | 73.85        | 68.1         |

Table 3.1: Baseline results of models trained on single languages from UD v1.3. Our models use only the orthographic representation of tokenized words during inference and works without a separate POS tagger.

In the following sections we will show impact of different network parameters on the result.

### 3.1.1 Impact of *reader* and *predictor*

We first inspect the impact of different *reader* and *predictor* settings.

| training inputs                   | Czech       |             | English     |             | Polish      |             |
|-----------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                                   | UAS         | LAS         | UAS         | LAS         | UAS         | LAS         |
| Gold POS tags                     |             |             |             |             |             |             |
| base word                         | <b>91.7</b> | <b>88</b>   | <b>88.6</b> | 85.1        | <b>93.4</b> | <b>89.3</b> |
| Predicted POS tags or no POS tags |             |             |             |             |             |             |
| words                             | 82.4        | 72.1        | 81.9        | 74.7        | 74.6        | 61.6        |
| chars, soft att.                  | <b>90.1</b> | 85.7        | 86.5        | 82.1        | 89.1        | 82.5        |
| chars, tags, soft att.            | 89.6        | 82.8        | 86.2        | 81.3        | 90.4        | 83.9        |
| chars, tags, hard att.            | <b>90.1</b> | <b>86.7</b> | <b>87.6</b> | <b>83.6</b> | <b>91.3</b> | <b>86</b>   |

Table 3.2: Model performance on selected languages and different training inputs. Experiments were run on **UD v 1.2**

In the upper part "Gold POS tags" we show a favorable situation when we have access to gold pos tag (approved by a human). The inputs are base words embedded in a vector (we replace the *reader* subnetwork with simple, learnable table lookup). This setting allowed us to obtain best results (especially for Polish) but gold pos tag are not available in a real-word setting so this result is noted just as a upper band of the model performance.

The "words" and "chars, soft att." settings don't use a *predictor*. Using just words as inputs (embedded in the same manner as base word) gives worst results. This is easy to predict because most words in the dataset occurs only once and so the network doesn't have a clue about role of most of the words. Instead, using character-level input allows us to obtain much better results, especially because in morphosyntactically rich languages the word spelling catches many word features that can be used in deducing the pos tag information.

The result of adding a *predictor* depends on the attention type used. Used with soft attention it decreased the results in Czech and English, but increased in Polish. When hard attention is used it can obtain the best results in all cases. The intuition behind these results is that with hard attention the labeller only sees annotation of the chosen head, and the error is backpropagated only to this element. On the other hand with soft attention the head annotation refers to many incorrect words introducing more noise. Additionally the soft attention labeller error is backpropagated through all words, which can introduce even more noise on the neural weights.

### 3.1.2 Impact of decoding algorithm

In the Table 3.3 we show the comparison of greedy and Chu-Liu-Edmonds decoding algorithms.

| language      | Greedy UAS   | Edmonds UAS |
|---------------|--------------|-------------|
| Czech         | <b>91.41</b> | 91.40       |
| Polish        | <b>90.26</b> | 90.17       |
| Russian       | <b>83.29</b> | 83.23       |
| German        | <b>82.67</b> | 82.66       |
| English       | <b>87.44</b> | 87.41       |
| French        | <b>87.25</b> | 87.18       |
| Ancient Greek | <b>78.96</b> | 78.92       |

Table 3.3: The comparison of greedy and edmonds decoding algorithms

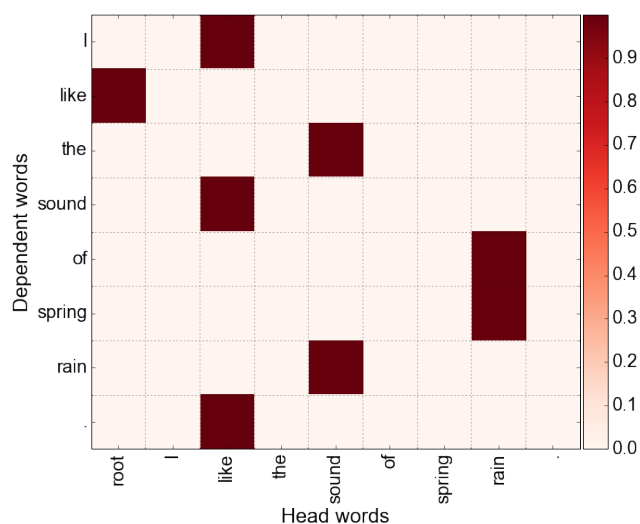


Figure 3.1: A typical scorer output showing probabilities assigned for head-word location.



| Language | % of recurrent state size | UAS          | LAS          |
|----------|---------------------------|--------------|--------------|
| Polish   | 100%                      | 90.26        | 85.32        |
|          | 50%                       | <b>90.54</b> | <b>85.64</b> |
|          | 25%                       | 89.07        | 82.80        |
| Czech    | 100%                      | <b>91.41</b> | <b>88.18</b> |
|          | 50%                       | 89.98        | 85.92        |
|          | 25%                       | 85.56        | 79.16        |

Table 3.4: Impact of the *tagger* recurrent state size on the performance. We report this value as percent of baseline rnn size.

depict each word as combination of those units. Of course to be able to present every word, each single letter also has to be treated as small subword unit. For example lets consider two polish words: *lepszy* and *najlepszy* (better, the best). The *naj* prefix in polish used with adjective means that it is in superlative degree, so we can imagine that this letter composition is pretty common and thus will be treated as subword unit. Then those two words will be shown to the network as *naj-l-e-p-s-z-y* and *l-e-p-s-z-y* (- acting as subword delimiter).

Our word pieces model are obtained as follows: first we compute word frequency dictionary and initialize pieces dictionary with all characters in the dataset. From frequency dictionary we take the most common character bigram and we treat it as new character and replace all previous occurrences with this new unit. The number of such iterations is a parameter given by the user. To convert word to pieces we use a greedy strategy. For word  $w$  we find the longest piece  $p$  that is prefix of  $w$ . The next step takes  $w$  suffix of size  $\|w\| - \|p\|$ .

| #pieces | UAS          | LAS          |
|---------|--------------|--------------|
| Polish  |              |              |
| 25      | 90.29        | 84.91        |
| 50      | <b>90.40</b> | <b>85.46</b> |
| 75      | 90.23        | 84.87        |
| 100     | 89.97        | 84.44        |
| Czech   |              |              |
| 25      | 90.29        | 86.50        |
| 50      | 90.03        | 86.08        |
| 75      | 90.17        | 86.40        |
| 100     | <b>90.84</b> | <b>87.31</b> |

Figure 3.3: Results on word pieces model. #pieces denotes how many new multi-character tokens were used. To convert word to pieces we use a greedy algorithm in which we choose the longest piece that is equal to prefix of word.

In Table 3.3 we have depicted the results of word pieces model for different number of additional multi-character tokens. There is only minor improvement for Polish while for Czech we have actually worse results. Additionally this input model is not well suited for multilanguage parser, because the pieces frequency is different for each language.

### 3.2 Multilanguage training

The main difference between multilanguage and single language configurations is that we have to decide which parts of the network should have shared weights. We tested it on 3 language pairs: Polish-Czech, Russian-Czech, Polish-Russian. They all are from the same Slavic family, but Russian uses Cyrillic alphabet as opposed to Latin script for Polish and Czech. The main advantage of Czech language that made us choose it as auxiliary language is the quality and quantity of training data - it is over 6 times bigger than Polish and Russian datasets combined. The results of different sharing strategies are depicted in Table 3.2.

| Shared parts                                 | Main lang | Aux lang | UAS          | LAS          |
|--|-----------|----------|--------------|--------------|
| -  | Polish    | -        | 90.31        | 85.21        |
| <i>Parser</i>                                | Polish    | Czech    | 90.72        | 85.57        |
| <i>Tagger, Parser</i>                        | Polish    | Czech    | 91.19        | 86.37        |
| <i>Tagger, POS Predictor, Parser</i>         | Polish    | Czech    | 91.65        | 86.88        |
| <i>Reader, Tagger, POS Predictor, Parser</i> | Polish    | Czech    | <b>91.91</b> | <b>87.77</b> |
| <i>Parser</i>                                | Polish    | Russian  | 90.31        | 85.07        |
| <i>Tagger, POS Predictor, Parser</i>         | Polish    | Russian  | <b>91.34</b> | <b>86.36</b> |
| <i>Reader, Tagger, POS Predictor, Parser</i> | Polish    | Russian  | 89.16        | 82.94        |
| -  | Russian   | -        | 83.43        | 79.24        |
| <i>Parser</i>                                | Russian   | Czech    | 83.15        | 78.69        |
| <i>Tagger, POS Predictor, Parser</i>         | Russian   | Czech    | 83.91        | 79.79        |
| <i>Reader, Tagger, POS Predictor, Parser</i> | Russian   | Czech    | <b>84.78</b> | <b>80.35</b> |

Table 3.5: Impact of parameter sharing strategies on main language parsing accuracy when multilingual training is used for additional supervision.

Multilingual training (Table 3.2) improves the performance on low-resource languages. We observe that the optimal amount of parameter sharing depends on the similarity between languages and corpus size – while it is beneficial to share all parameters of the PL-CZ and RU-CZ parser, the PL-RU parser works best if the reader subnetworks are separated.

Similarly as in the single language configuration we have tested the impact of *tagger* recurrent state size on the results (see Table 3.2).

| % of recurrent state size | UAS          | LAS          |
|---------------------------|--------------|--------------|
| 200%                      | 91.45        | 86.36        |
| 100%                      | <b>91.65</b> | <b>86.88</b> |
| 50%                       | 89.53        | 83.66        |
| 25%                       | 87.28        | 78.93        |

Table 3.6: Impact of the *tagger* recurrent state size on the performance of multilanguage pl-cz model. We report this value as percent of baseline rnn size.



The results are computed on best Polish-Czech model (with share all strategy). We can see that the recurrent state size have similar behaviour as in the single language Czech model - the 50% size has significant drop in the performance. It is worth noting that increasing the size changes the result only slightly (and that they are actually worse). This mean that chosen *tagger* size is well-suited for both single and multi language parsers.

### 3.2.1 Polish-Polish multilang parser

After investigating results of the parser (this problem applies to both single and multi parser) we have found that if the *tagger* assigns wrong POS tag to a particular word it can mislead the scorer and labeller producing wrong parse tree. In many cases after replacing wrong word with its synonym such that the *tagger* will assign the right POS tag - the tree becomes a valid parse tree. For an example see Figure 3.4.

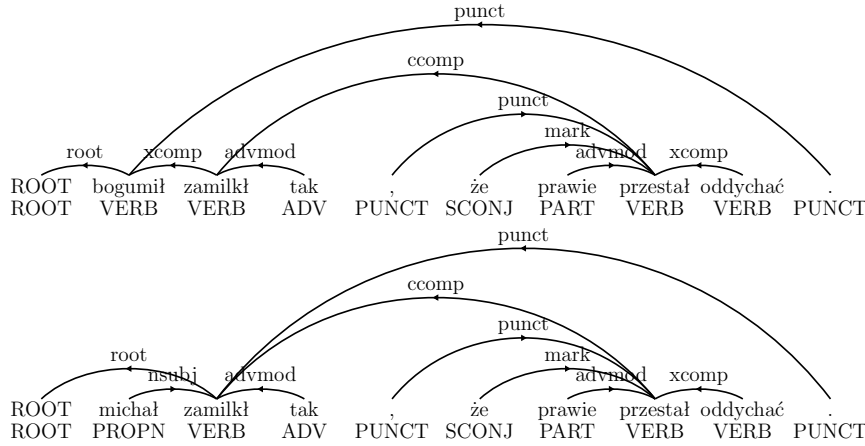


Figure 3.4: In the upper panel the sentence "Bogumił felt silent so much that he almost stopped breathing." is incorrectly parsed. We attribute the error to the *reader* or to the *tagger* because the predicted POS tag for "Bogumił", a male given name, is wrong (VERB instead of NOUN). In the lower panel we replace "Bogumił" with the more common "Michał" the network assigns a proper tag and the whole sentence becomes a valid parse tree.

To prevent this behaviour on Polish language we have decided to use additional data for training the *predictor* subnetwork. We used the *milionowy podkorpus* (million subcorpus) of the National Corpus of Polish Language [Przepiórkowski et al., 2012]. We converted this data to UD conllu format by filling the dependency information with a blank values and transforming the polish-style POS tags to UD-style using conversion table provided by UD authors. Next we used our multilingual training procedure in which the first language was a standard UD Polish dataset, and the second one was our converted polish subcorpus. The difference with previous attempts is that we have to ignore loss on dependency head  $L_h$  and dependency label  $L_l$  for the second language.

This model allowed us to obtain our best results: **92.53% UAS** and **88.84% LAS**.



## Chapter 4

# Summary

We have demonstrated a nonprojective graph-based dependency parser implemented as a single deep neural network that directly produces parse trees from characters and does not require other NLP tools such as a POS tagger. The parser is trained in an end-to-end manner, and has separate cost terms that pertain to label accuracy, head word localization and optionally POS tagging. On morphologically rich languages the parser is competitive with other state-of-the-art parsers.

With additional experiments we have established that the multilingual training performance depends heavily on degree of parameter sharing, and can differ depending on language similarity and corpus size. The decoding algorithm used does not have a big impact on the performance of the system, so simple greedy one is sufficient.

By training part of the network with additional, language-specific data we were able to exceed the state-of-the-art results by a margin of over 0.5 percent point.



# Bibliography

- [Alberti et al., 2017] Alberti, C. et al. (2017). SyntaxNet Models for the CoNLL 2017 Shared Task. *arXiv:1703.04929*.
- [Andor et al., 2016] Andor, D., Alberti, C., Weiss, D., Severyn, A., Presta, A., Ganchev, K., Petrov, S., and Collins, M. (2016). Globally Normalized Transition-Based Neural Networks. *arXiv:1603.06042 [cs]*. 00001 arXiv: 1603.06042.
- [Bahdanau et al., 2014] Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv:1409.0473 [cs, stat]*. 00000 arXiv: 1409.0473.
- [Bejček et al., 2013] Bejček, E., Hajičová, E., Hajič, J., Jínová, P., Kettnerová, V., Kolářová, V., Mikulová, M., Mírovský, J., Nedoluzhko, A., Panevová, J., Poláková, L., Ševčíková, M., Štěpánek, J., and Zikánová, Š. (2013). Prague dependency treebank 3.0.
- [Caruana, 1997] Caruana, R. (1997). Multitask learning. *Machine Learning*, 28(1):41–75.
- [Chen and Manning, 2014] Chen, D. and Manning, C. D. (2014). A Fast and Accurate Dependency Parser using Neural Networks. In *EMNLP*, pages 740–750. 00188.
- [Cho et al., 2014] Cho, K., van Merriënboer, B., Gülçehre, Ç., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078.
- [Chorowski et al., 2014] Chorowski, J., Bahdanau, D., Cho, K., and Bengio, Y. (2014). End-to-end Continuous Speech Recognition using Attention-based Recurrent NN: First Results. *arXiv:1412.1602 [cs, stat]*. 00000 arXiv: 1412.1602.
- [Covington, 2001] Covington, M. A. (2001). A fundamental algorithm for dependency parsing. In *Proceedings of the 39th annual ACM southeast conference*, pages 95–102. Citeseer. 00206.
- [De Marneffe et al., 2006] De Marneffe, M.-C., MacCartney, B., Manning, C. D., et al. (2006). Generating typed dependency parses from phrase structure parses. In *Proceedings of LREC*, volume 6, pages 449–454. Genoa.

- [Edmonds, 1966] Edmonds, J. (1966). Optimim Branchings. *JOURNAL OF RESEARCH of the National Bureau of Standards - B.*, 71B(4):233–240. 00000.
- [Gabow et al., 1986] Gabow, H. N., Galil, Z., Spencer, T., and Tarjan, R. E. (1986). Efficient algorithms for finding minimum spanning trees in undirected and directed graphs. *Combinatorica*, 6(2):109–122.
- [Goodfellow et al., 2013] Goodfellow, I., Warde-Farley, D., Mirza, M., Courville, A., and Bengio, Y. (2013). Maxout Networks. In *ICML*, pages 1319–1327. 00106.
- [Hochreiter and Schmidhuber, 1997] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Comput.*, 9(8):1735–1780.
- [Kim et al., 2015] Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2015). Character-aware neural language models. *arXiv preprint arXiv:1508.06615*. 00011  
bibtex: kim\_character\_2015a.
- [Nair and Hinton, 2010] Nair, V. and Hinton, G. E. (2010). Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814.
- [Nivre, 2008] Nivre, J. (2008). Algorithms for Deterministic Incremental Dependency Parsing. *Comput. Linguist.*, 34(4):513–553. 00220.
- [Nivre et al., 2015] Nivre, J. et al. (2015). Universal Dependencies 1.2. <http://universaldependencies.github.io/docs/>.
- [Nivre et al., 2005] Nivre, J., Hall, J., Nilsson, J., Chanev, A., Eryigit, G., Kübler, S., Marinov, S., and Marsi, E. (2005). MaltParser: A language-independent system for data-driven dependency parsing. *Natural Language Engineering*, page 1. 00666.
- [Petrov et al., 2011] Petrov, S., Das, D., and McDonald, R. (2011). A universal part-of-speech tagset. In *LREC 2012*.
- [Przepiórkowski et al., 2012] Przepiórkowski, A., Bańko, M., Górski, R. L., and Lewandowska-Tomaszczyk, B., editors (2012). *Narodowy Korpus Języka Polskiego*. Wydawnictwo Naukowe PWN, Warsaw.
- [Schuster and Paliwal, 1997] Schuster, M. and Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681. 00157.
- [Sennrich et al., 2015] Sennrich, R., Haddow, B., and Birch, A. (2015). Neural Machine Translation of Rare Words with Subword Units. *ArXiv e-prints*.
- [Snoek et al., 2012] Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959. 00414.

- [Srivastava et al., 2014] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15:1929–1958. 00282.
- [Straka et al., 2015] Straka, M., Hajič, J., Straková, J., and Jan Hajič, j. (2015). Parsing universal dependency treebanks using neural networks and search-based oracle. In *14th International Workshop on Treebanks and Linguistic Theories (TLT 2015)*, pages 208–220, Warszawa, Poland. IPIPAN, IPIPAN.
- [Szegedy et al., 2014] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S. E., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2014). Going deeper with convolutions. *CoRR*, abs/1409.4842.
- [Szegedy et al., 2015] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2015). Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567.
- [Vinyals et al., 2015] Vinyals, O., Fortunato, M., and Jaitly, N. (2015). Pointer networks. In *Advances in Neural Information Processing Systems*, pages 2674–2682. 00010.
- [Xu et al., 2015] Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R., and Bengio, Y. (2015). Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. *arXiv:1502.03044 [cs]*. 00001 arXiv: 1502.03044.
- [Zeiler, 2012] Zeiler, M. D. (2012). ADADELTA: An Adaptive Learning Rate Method. *arXiv:1212.5701 [cs]*. 00017 arXiv: 1212.5701.
- [Zeman, 2008] Zeman, D. (2008). Reusable tagset conversion using tagset drivers. In *LREC*.