

# SEMI-SUPERVISED AUTOENCODING PROJECTIVE DEPENDENCY PARSING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We describe two end-to-end autoencoding models for semi-supervised graph-based dependency parsing. The first model is a Local Autoencoding Parser (LAP) encoding the input using continuous latent variables in a sequential manner; The second model is a Global Autoencoding Parser (GAP) encoding the input into dependency trees as latent variables, with exact inference. Both models consist of two parts: an encoder enhanced by deep neural networks (DNN) that can utilize the contextual information to encode the input into latent variables, and a decoder which is a generative model able to reconstruct the input. Both LAP and GAP admit a unified structure with different loss functions for labeled and unlabeled data with shared parameters. We conducted experiments on WSJ and UD dependency parsing data sets, showing that our models can exploit the unlabeled data to boost the performance given a limited amount of labeled data.

## 1 INTRODUCTION

Dependency parsing captures bi-lexical relationships by constructing directional arcs between words, defining a head-modifier syntactic structure for sentences, as shown in Figure 1. Dependency trees are fundamental for many downstream tasks such as semantic parsing (Reddy et al., 2016; Marcheggiani & Titov, 2017), machine translation (Bastings et al., 2017; Ding & Palmer, 2007), information extraction (Culotta & Sorensen, 2004; Liu et al., 2015) and question answering (Cui et al., 2005). As a result, efficient parsers (Kiperwasser & Goldberg, 2016; Dozat & Manning, 2017; Dozat et al., 2017; Ma et al., 2018) have been developed using various neural architectures.

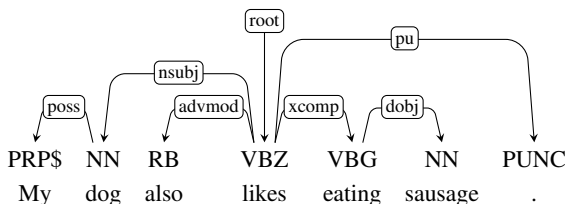


Figure 1: A dependency tree: directional arcs represent head-modifier relation between words.

While supervised approaches have been very successful, they require large amounts of labeled data, particularly when neural architectures are used. Syntactic annotation is notoriously difficult and requires specialized linguistic expertise, posing a serious challenge for low-resource languages. Semi-supervised parsing aims to alleviate this problem by combining a small amount of labeled data and a large amount of unlabeled data, to improve parsing performance over labeled data alone. Traditional semi-supervised parsers use unlabeled data to generate additional features, assisting the learning process (Koo et al., 2008), together with different variants of self-training (Søgaard & Rishøj, 2010). However, these approaches are usually pipe-lined and error-propagation may occur.

In this paper, we propose two end-to-end semi-supervised parsers based on probabilistic autoencoder models illustrated in Figure 3, Locally Autoencoding Parser (LAP) and Globally Autoencoding Parser (GAP). In LAP, continuous latent variables are used to support tree inference by providing a better representation, while in GAP, the latent information forms a probability distribution over

dependency trees corresponding to the input sentence. A similar idea has been proposed by Corro & Titov (2018), but our GAP model differs fundamentally from their parser, as GAP does not sample from the posterior of the latent tree structure to approximate the Evidence Lower Bound (ELBO). Instead it relies on a tractable algorithm to directly compute the posterior to calculate the ELBO.

We summarize our contributions as follows:

1. We proposed two autoencoding parsers for semi-supervised dependency parsing, with complementary strengths, trading off speed vs. accuracy;
2. We propose a tractable inference algorithm to compute the expectation and marginalization of the latent dependency tree posterior analytically for GAP, avoiding sampling from the posterior to approximate the expectation (Corro & Titov, 2018);
3. We show improved performance of both LAP and GAP with unlabeled data on WSJ and UD data sets empirically, and improved results of GAP comparing to a recently proposed semi-supervised parser (Corro & Titov, 2018).

## 2 RELATED WORK

Most dependency parsing studies fall into two major groups: graph-based and transition-based (Kubler et al., 2009). Graph-based parsers (McDonald, 2006) regard parsing as a structured prediction problem to find the most probable tree, while transition-based parsers (Nivre, 2004; 2008) treat parsing as a sequence of actions at different stages leading to a dependency tree.

While earlier works relied on manual feature engineering, in recent years the hand-crafted features were replaced by embeddings and deep neural architectures, leading to improved performance in both graph-based parsing (Nivre, 2014; Pei et al., 2015) and transition-based parsing (Chen & Manning, 2014; Dyer et al., 2015; Weiss et al., 2015). More recent works rely on neural architectures for learning a representation for scoring structural decisions Andor et al. (2016); Kiperwasser & Goldberg (2016); Wiseman & Rush (2016).

The annotation difficulty for this task, has also motivated work on unsupervised (grammar induction) and semi-supervised approaches to parsing (Tu & Honavar, 2012; Jiang et al., 2016; Koo et al., 2008; Li et al., 2014; Kiperwasser & Goldberg, 2015; Cai et al., 2017; Corro & Titov, 2018).

Similar to other structured prediction tasks, directly optimizing the objective is difficult when the underlying probabilistic model requires marginalizing over the dependency trees. Variational approaches are a natural way for alleviating this problem, as they try to improve the lower bound of the original objective, and were applied in several recent NLP works (Stratos, 2019; Kim et al., 2019b; Chen et al., 2018; Kim et al., 2019b;a). Variational Autoencoder (VAE) (Kingma & Welling, 2014) is particularly useful for latent representation learning, and is studied in semi-supervised context as the Conditional VAE (CVAE) (Sohn et al., 2015).

The work mostly related to ours is (Corro & Titov, 2018) as they consider the dependency tree as the latent variable, but their work takes a second approximation to the variational lower bound by an extra step to sample from the latent dependency tree, without identifying a tractable inference. We show that with the given structure, exact inference on the lower bound is achievable without approximation by sampling, which tightens the lower bound.

## 3 GRAPH-BASED DEPENDENCY PARSING

A dependency graph of a sentence can be regarded as a directed tree spanning all the words of the sentence, including a special “word”—the ROOT—to originate out. Assuming a sentence length of  $l$ , a dependency tree can be denoted as  $\mathcal{T} = (\langle h_1, m_1 \rangle, \dots, \langle h_{l-1}, m_{l-1} \rangle)$ , where  $h_t$  is the index in the sequence of the head word of the dependency connecting the  $t$ th word  $m_t$  as a modifier.

Our graph-based parser is constructed by following the standard structured prediction paradigm (McDonald et al., 2005; Taskar et al., 2005). In inference, based on the parameterized scoring function  $\mathcal{S}_\Lambda$  with parameter  $\Lambda$ , the parsing problem is formulated as finding the most probable directed spanning tree for a given sentence  $x$ :

$$\mathcal{T}^* = \arg \max_{\tilde{\mathcal{T}} \in \mathbb{T}} \mathcal{S}_\Lambda(x, \tilde{\mathcal{T}}),$$

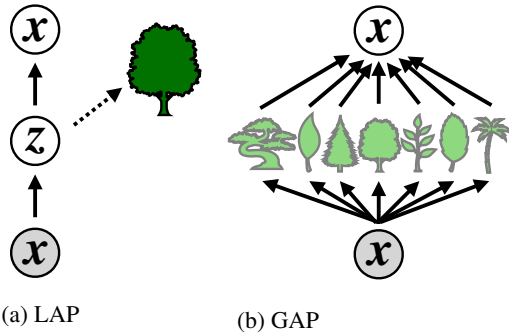


Figure 3: Illustration of two different parsers. (a) LAP uses continuous latent variable to form the dependency tree (b) GAP treats the dependency tree as the latent variable.

	Root	$x_1$	...	$x_t$	...	$x_{l-1}$	
Root	0	The score of the (t, t-1) right arc					
$x_1$		0					
$\vdots$			0				
$x_t$				0	$s_{(t,t-1)}$		
$\vdots$						0	
$x_{l-1}$						0	
	The score of the (t-1, t) left arc						0

Figure 4: In this illustration of the arc scoring matrix, each entry represents the  $(h(\text{head}) \rightarrow m(\text{modifier}))$  score.

where  $\mathcal{T}^*$  is the highest scoring parse tree and  $\mathbb{T}$  is the set of all valid trees for the sentence  $\mathbf{x}$ .

It is common to factorize the score of the entire graph into the summation of its substructures: the individual arc scores (McDonald et al., 2005):

$$S_{\Lambda}(\mathbf{x}, \tilde{\mathcal{T}}) = \sum_{(h,m) \in \tilde{\mathcal{T}}} s_{\Lambda}(h, m) = \sum_{t=1}^l s_{\Lambda}(h_t, m_t),$$

where  $\tilde{\mathcal{T}}$  represents the candidate parse tree, and  $s_{\Lambda}$  is a function scoring each individual arc.  $s_{\Lambda}(h, m)$  describes the likelihood of forming an arc from the head  $h$  to its modifier  $m$  in the tree. Through out this paper, the scoring is based on individual arcs, as we focus on *first order* parsing.

### 3.1 SCORING FUNCTION USING NEURAL ARCHITECTURE

We used the same neural architecture as that in Kiperwasser & Goldberg (2016)’s study. We first use a bi-LSTM model to take as input  $\mathbf{u}_t = [p_t; e_t]$  at position  $t$  to incorporate contextual information, by feeding the word embedding  $e_t$  concatenated with the POS (part of speech) tag embeddings  $p_t$  of each word. The bi-LSTM then projects  $\mathbf{u}_t$  as  $\mathbf{o}_t$ .

Subsequently a nonlinear transformation is applied on these projections. Suppose the hidden states generated by the bi-LSTM are  $[\mathbf{o}_{root}, \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_t, \dots, \mathbf{o}_l]$ , for a sentence of length  $l$ , we compute the arc scores by introducing parameters  $\mathbf{W}_h, \mathbf{W}_m, \mathbf{w}$  and  $\mathbf{b}$ , and transform them as follows:

$$\mathbf{r}_t^{h-arc} = \mathbf{W}_h \mathbf{o}_t; \mathbf{r}_t^{m-arc} = \mathbf{W}_m \mathbf{o}_t,$$

$$s_{\Lambda}(h, m) = \mathbf{w}^T (\tanh(\mathbf{r}_h^{h-arc} + \mathbf{r}_m^{m-arc} + \mathbf{b})).$$

In this formulation, we first use two parameters to extract two different representations that carry two different types of information: a head seeking for its modifier (h-arc); as well as a modifier seeking for its head (m-arc). Then a nonlinear function maps them to an arc score.

For a single sentence, we can form a scoring matrix as shown in Figure 4, by filling each entry in the matrix using the score we obtained. Therefore, the scoring matrix is used to represent the head-modifier arc score of all the possible arcs connecting words in a sentence (Zheng, 2017).

Using the scoring arc matrix, we build graph-based parsers. Since exploring neural architectures for scoring is not our focus, we did not explore other architectures, however performance shall be further improved using advanced neural architectures (Dozat & Manning, 2017; Dozat et al., 2017).

## 4 PRELIMINARIES: VARIATIONAL AUTO-ENCODER AND TREE CRF

**Variational Autoencoder (VAE).** The typical VAE is a directed graphical model with Gaussian latent variables, denoted by  $\mathbf{z}$ . A generative process first generates a set of  $\mathbf{z}$  from the prior distribution  $\pi(\mathbf{z})$  and the data  $\mathbf{x}$  is generated as  $P_{\theta}(\mathbf{x}|\mathbf{z})$  parameterized by  $\theta$  given input  $\mathbf{x}$ . In our scenario,  $\mathbf{x}$  is an input sequence and  $\mathbf{z}$  is a sequence of latent variables corresponding to it.

The VAE framework seeks to maximize the complete log-likelihood  $\log P(\mathbf{x})$  by marginalizing out the latent variable  $\mathbf{z}$ . Since direct parameter estimation of  $\log P(\mathbf{x})$  is usually intractable, a common solution is to maximize its Evidence Lower Bound (ELBO) by introducing an auxiliary posterior  $Q(\mathbf{x}|\mathbf{z})$  distribution that encodes the input into the latent space.

**Tree Conditional Random Field.** Linear chain CRF models an input sequence  $\mathbf{x} = (x_1 \dots x_l)$  of length  $l$  with labels  $\mathbf{y} = (y_1 \dots y_l)$  with globally normalized probability

$$P(\mathbf{y}|\mathbf{x}) = \frac{\exp \mathcal{S}(\mathbf{x}, \mathbf{y})}{\sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} \exp \mathcal{S}(\mathbf{x}, \tilde{\mathbf{y}})},$$

where  $\mathcal{Y}$  is the set of all the possible label sequences, and  $\mathcal{S}(\mathbf{x}, \mathbf{y})$  the scoring function, usually decomposed as emission ( $\sum_{i=1}^l s(x_i, y_i)$ ) and transition ( $\sum_{i=1}^l s(y_i, y_{i+1})$ ) for *first order* models.

Tree CRF models generalize linear chain CRF to trees. For dependency trees, if POS tags are given, the tree CRF model tries to resolve which node pairs should be connected with direction, such that the arcs form a tree. The potentials in the dependency tree take an exponential form, thus the conditional probability of a parse tree  $\mathcal{T}$ , given the sequence, can be denoted as:

$$P(\mathcal{T}|\mathbf{x}) = \frac{\exp \mathcal{S}(\mathbf{x}, \mathcal{T})}{Z(\mathbf{x})}, \quad (1)$$

where  $Z(\mathbf{x}) = \sum_{\tilde{\mathcal{T}} \in \mathbb{T}(\mathbf{x})} \exp \mathcal{S}(\mathbf{x}, \tilde{\mathcal{T}})$  is the partition function that sums over all possible valid dependency trees in the set  $\mathbb{T}(\mathbf{x})$  of the given sentence  $\mathbf{x}$ .

## 5 LOCALLY AUTOENCODING PARSER (LAP)

We extend the original VAE model for sequence labeling (Chen et al., 2018) to dependency parsing by building a latent representation position-wise to form a sequential latent representation.

It has been shown that under the VAE framework the latent representation can reflect the desired properties of the raw input (Kingma & Welling, 2014). This inspired us to use the continuous latent variable as neural representations for the dependency parsing task. Typically, each token in the sentence is represented by its latent variable  $z_t$ , which is a high-dimensional Gaussian variable. This configuration on the one hand ensures the continuous latent variable retains the contextual information from lower-level neural models to assist finding its head or its modifier; on the other hand, it forces tokens of similar properties closer in the euclidean space. We adjust the original VAE setup in our semi-supervised task by considering examples with labels, similar to recent conditional variational formulations (Sohn et al., 2015; Miao & Blunsom, 2016; Zhou & Neubig, 2017).

We propose a full probabilistic model for any certain sentence  $\mathbf{x}$ , with the unified objective to maximize for supervised and unsupervised parsing as follows:

$$\mathcal{J} = \log P_{\theta}(\mathbf{x}) P_{\omega}^{\epsilon}(\mathcal{T}|\mathbf{x}), \quad \epsilon = \begin{cases} 1, & \text{if } \mathcal{T} \text{ exists,} \\ 0, & \text{otherwise.} \end{cases}$$

This objective can be interpreted as follows: if the training example has a golden tree  $\mathcal{T}$  with it, then the objective is the log joint probability  $P_{\theta, \omega}(\mathcal{T}, \mathbf{x})$ ; if the golden tree is missing, then the objective is the log marginal probability  $P_{\theta}(\mathbf{x})$ . The probability of a certain tree is modeled by a tree-CRF in Eq. 1 with parameters  $\omega$  as  $P_{\omega}(\mathcal{T}|\mathbf{x})$ . Given the assumed generative process  $P(\mathbf{x}|\mathbf{z})$ , directly optimizing this objective is intractable, we instead optimize its ELBO (We show the details in the appendix, proving  $\mathcal{J}_{lap}$  is the ELBO of  $\mathcal{J}$  in Lemma A.1):

$$\mathcal{J}_{lap} = \mathbb{E}_{\mathbf{z} \sim Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log P_{\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{KL}(Q_{\phi}(\mathbf{z}|\mathbf{x})||P_{\theta}(\mathbf{z})) + \epsilon \mathbb{E}_{\mathbf{z} \sim Q_{\phi}(\mathbf{z}|\mathbf{x})} [\log P_{\omega}(\mathcal{T}|\mathbf{z})].$$

## 6 GLOBALLY AUTOENCODING PARSER (GAP)

Instead of autoencoding the input locally at the sequence level, we could alternatively directly regard the dependency tree as the structured latent variable to reconstruct the input sentence, by building a model containing both a discriminative component and a generative component. The discriminative component builds a neural CRF model for dependency tree construction, and the generative model reconstructs the sentence from the factor graph as a Bayesian network, by assuming a generative

process in which each head generates its modifier. Concretely, the latent variable in this model is the dependency tree structure.

### 6.1 DISCRIMINATIVE COMPONENT: THE ENCODER

We model the discriminative component in our model as  $P_{\Phi}(\mathcal{T}|\mathbf{x})$  parameterized by  $\Phi$ , taking the same form as in Eq. 1. Typically in our model,  $\Phi$  are the parameters of the underlying neural networks, whose architecture is described in Sec. 3.1.

### 6.2 GENERATIVE COMPONENT: THE DECODER

We use a set of conditional categorical distributions to construct our Bayesian network decoder. More specifically, using the head  $h$  and modifier  $m$  notation, each head reconstructs its modifier with the probability  $P(m_t|h_t)$  for the  $t$ th word in the sentence (0th word is always the special ‘‘ROOT’’ word), which is parameterized by the set of parameters  $\Theta$ . Given  $\Theta$  as a matrix of  $|\mathcal{V}|$  by  $|\mathcal{V}|$ , where  $|\mathcal{V}|$  is the vocabulary size,  $\theta_{mh}$  is the item on row  $m$  column  $h$  denoting the probability that the head word  $h$  would generate  $m$ . In addition, we have a simplex constraint  $\sum_{m \in \mathcal{V}} \theta_{mh} = 1$ . The probability of reconstructing the input  $\mathbf{x}$  as modifiers  $\mathbf{m}$  in the generative process is

$$P_{\Theta}(\mathbf{m}|\mathcal{T}) = \prod_t P(m_t|h_t) = \prod_t \theta_{m_t h_t},$$

where  $l$  is the sentence length and  $P(m_t|h_t)$  represents the probability a head generating its modifier.

### 6.3 A UNIFIED SUPERVISED AND UNSUPERVISED LEARNING FRAMEWORK

With the design of the discriminative component and the generative component of the proposed model, we have a unified learning framework for sentences with or without golden parse tree.

The complete data likelihood of a given sentence, if the golden tree is given, is

$$\begin{aligned} P_{\Theta, \Phi}(\mathbf{m}, \mathcal{T}|\mathbf{x}) &= P_{\Theta}(\mathbf{m}|\mathcal{T})P_{\Phi}(\mathcal{T}|\mathbf{x}) \\ &= \left[ \prod_{t=1}^l P(m_t|h_t) \right] \frac{\exp \mathcal{S}_{\Phi}(\mathbf{x}, \mathcal{T})}{Z(\mathbf{x})} \\ &= \frac{\exp \sum_{(h,m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m)}{Z(\mathbf{x})}, \end{aligned}$$

where  $s'_{\Phi, \Theta}(h, m) = s_{\Phi}(h, m) + \log \theta_{mh}$ , with  $\mathbf{m}, \mathbf{x}$  and  $\mathcal{T}$  all observable.

For unlabeled sentences, the complete data likelihood can be obtained by marginalizing over all the possible parse trees in the set  $\mathbb{T}(\mathbf{x})$ :

$$\begin{aligned} P_{\Theta, \Phi}(\mathbf{m}|\mathbf{x}) &= \sum_{\mathcal{T} \in \mathbb{T}(\mathbf{x})} P_{\Theta, \Phi}(\mathbf{m}, \mathcal{T}|\mathbf{x}) \\ &= \frac{U(\mathbf{x})}{Z(\mathbf{x})}, \end{aligned}$$

where  $U(\mathbf{x}) = \sum_{\mathcal{T} \in \mathbb{T}(\mathbf{x})} \exp \sum_{(h,m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m)$ .

We adapted a variant of Eisner (1996)’s algorithm to marginalize over all possible trees to compute both  $Z$  and  $U$ , as  $U$  has the same structure as  $Z$ , assuming a projective tree.

We use log-likelihood as our objective function. The objective for a sentence with golden tree is:

$$\begin{aligned} \mathcal{J}_l &= \log P_{\Theta, \Phi}(\mathbf{m}, \mathcal{T}|\mathbf{x}) \\ &= \sum_{(h,m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m) - \log Z(\mathbf{x}) \end{aligned}$$

**Algorithm 1** Learning Algorithm for GAP

- 
- 1: Initialize the parameter  $\Theta$  in the decoder with the labeled data set  $\{\mathbf{x}, \mathcal{T}\}^l$ .
  - 2: Initialize  $\Lambda$  in the encoder randomly.
  - 3: **for**  $t$  in *epochs* **do**
  - 4:     **for** sentence  $\mathbf{x}_i^l$  with golden parse tree  $\mathcal{T}_i^l$  in the labeled data set  $\{\mathbf{x}, \mathcal{T}\}^l$  **do**
  - 5:         Stochastically update the parameter  $\Lambda$  in the encoder using Adam while fixing the decoder.
  - 6:     **end for**
  - 7:     Initialize a Counting Buffer  $\mathcal{B}$
  - 8:     **for** unlabeled sentence  $\mathbf{x}_i^u$  in the unlabeled data set  $\{\mathbf{x}\}^u$  **do**
  - 9:         Compute the posterior  $Q(\mathcal{T})$  in an arc factored manner for  $\mathbf{x}_i^u$  tractably.
  - 10:         Compute the expectation of all possible  $(h(\text{head}) \rightarrow m(\text{modifier}))$  occurrence in the sentence  $\mathbf{x}$  based on  $Q(\mathcal{T})$ .
  - 11:         Update buffer  $\mathcal{B}$  using the expectation to the power for  $\frac{1}{1-\sigma}$  of all possible  $(h \rightarrow m)$ .
  - 12:     **end for**
  - 13:     Obtain  $\Theta$  globally and analytically based on the buffer  $\mathcal{B}$  and renew the decoder.
  - 14: **end for**
- 

If the input sentence does not have an annotated golden tree, then the objective is:

$$\begin{aligned} \mathcal{J}_u &= \log P_{\Theta, \Phi}(\mathbf{m}|\mathbf{x}) \\ &= \log U(\mathbf{x}) - \log Z(\mathbf{x}). \end{aligned} \quad (2)$$

Thus, during training, the objective function with shared parameters is chosen based on whether the sentence in the corpus has golden parse tree or not.

#### 6.4 LEARNING

Directly optimizing the loss in Eq.2 is difficult for the unlabeled data, and may lead to undesirable shallow local optima without any constraints. Instead, we derive the evidence lower bound (ELBO) of  $\log P_{\Theta, \Phi}(\mathbf{m}|\mathbf{x})$  as follows, by denoting  $Q(\mathcal{T}) = P_{\Theta, \Phi}(\mathcal{T}|\mathbf{m}, \mathbf{x})$  as the posterior:

$$\begin{aligned} \log P_{\Theta, \Phi}(\mathbf{m}|\mathbf{x}) &= \log \sum_{\mathcal{T}} Q(\mathcal{T}) \frac{P_{\Theta, \Phi}(\mathbf{m}, \mathcal{T}|\mathbf{x})}{Q(\mathcal{T})} \\ &= \log \mathbb{E}_{\mathcal{T} \sim Q(\mathcal{T})} \frac{P_{\Theta, \Phi}(\mathbf{m}, \mathcal{T}|\mathbf{x})}{Q(\mathcal{T})} \\ &\geq \mathbb{E}_{\mathcal{T} \sim Q(\mathcal{T})} \log \frac{P_{\Theta, \Phi}(\mathbf{m}, \mathcal{T}|\mathbf{x})}{Q(\mathcal{T})} \\ &= \mathbb{E}_{\mathcal{T} \sim Q(\mathcal{T})} [\log P_{\Theta}(\mathbf{m}|\mathcal{T})] - \mathbb{KL}[Q(\mathcal{T})||P_{\Phi}(\mathcal{T}|\mathbf{x})]. \end{aligned}$$

Instead of maximizing the log-likelihood directly, we alternatively maximize the ELBO, so our new objective function for unlabeled data becomes

$$\max_{\Theta, \Phi} \sum_i \mathbb{E}_{\mathcal{T} \sim Q(\mathcal{T})} [\log P_{\Theta}(\mathbf{m}|\mathcal{T})] - \mathbb{KL}[Q(\mathcal{T})||P_{\Phi}(\mathcal{T}|\mathbf{x})].$$

In addition, to account for the unambiguity in the posterior, we incorporate entropy regularization (Tu & Honavar, 2012) when applying our algorithm, by adding an entropy term  $-\sum_{\mathcal{T}} Q(\mathcal{T}) \log Q(\mathcal{T})$  with a non-negative factor  $\sigma$  when the input sentence does not have a golden tree. Adding this regularization term is equivalent as raising the expectation of  $Q(\mathcal{T})$  to the power of  $\frac{1}{1-\sigma}$ . We annealed  $\sigma$  from 1 to 0.3 from the beginning of training to the end, as in the beginning, the generative model is well initialized by sentences with golden trees that resolve disambiguity.

In practice, we found the model benefits more by fixing the parameter  $\Phi$  when the data is unlabeled and optimizing the ELBO w.r.t. the parameter  $\Theta$ . We attribute this to the strict convexity of the ELBO w.r.t.  $\Theta$ , by sketching the proof in the appendix. The details of training are shown in Alg. 1.

#### 6.5 TRACTABLE INFERENCE

The common approach to approximate the expectation of the latent variables from the posterior distribution  $Q(\mathcal{T})$  is via sampling in VAE-type models (Kingma & Welling, 2014). In a significant

contrast to that, we argue in this model the expectation of the latent variable (which is the dependency tree structure) is analytically tractable by designing a variant of the inside-outside algorithm (Eisner, 1996; Paskin, 2001) in an arc decomposed manner. We leave the detailed derivation in the appendix. A high-level explanation is that assuming the dependency tree is *projective*, specialized *belief propagation* algorithm exists to compute not only the *marginalization* but also the *expectation* analytically, making inference tractable.

## 7 EXPERIMENTS

### 7.1 EXPERIMENTAL SETTINGS

**Data sets** First we compared our models’ performance with strong baselines on the WSJ data set, which is the Stanford Dependency conversion (De Marneffe & Manning, 2008) of the Penn Treebank (Marcus et al., 1993) using the standard section split: 2-21 for training, 22 for development and 23 for testing. Second we evaluated our models on multiple languages, using data sets from UD (Universal Dependency) 2.3 (McDonald et al., 2013). Since semi-supervised learning is particularly useful for low-resource languages, we believe those languages in UD can benefit from our approach. The statistics of the data used in our experiments are described in Table 3 in appendix.

To simulate the low-resource language environment, we used 10% of the whole training set as the annotated, and the rest 90% as the unlabeled.

**Input Representation and Architecture** Since we use the same neural architecture in all of our models, we specify the details of the architecture once, as follows: The internal word embeddings have dimension 100 and the POS embeddings have dimension 25. The hidden layer of the bi-LSTM layer is of dimension 125. The nonlinear layers used to form the head and the modifier representation both have 100 dimension. For LAP, we use separate bi-LSTMs for words and POSs. In GAP, using “POS to POS” decoder only yield the satisfactory performance. This echoes the finding that complicated decoders may cause “posterior collapse” (van den Oord et al., 2017; Kim et al., 2018).

**Training** In the training phase, we use Adam (Kingma & Ba, 2014) to update all the parameters in both LAP and GAP, except the parameters in the decoder in GAP, which are updated by using their global optima in each epoch. We did not take efforts to tune models’ hyper-parameters and they remained the same across all the experiments.

### 7.2 SEMI-SUPERVISED DEPENDENCY PARSING ON WSJ DATA SET

We first evaluate our models on the WSJ data set and compared the model performance with other semi-supervised parsing models, including CRFAE (Cai et al., 2017), which is originally designed for dependency grammar induction but can be modified for semi-supervised parsing, and “differentiable Perturb-and-Parse” parser (DPPP) (Corro & Titov, 2018). To contextualize the results, we also experiment with the supervised neural margin-based parser (NMP) (Kiperwasser & Goldberg, 2016), neural tree-CRF parser (NTP) and the supervised version of LAP and GAP, with only the labeled data. To ensure a fair comparison, our experimental set up on the WSJ is identical as that in DPPP and we use the same 100 dimension skip-gram word embeddings employed in an earlier transition-based system (Dyer et al., 2015). We show our experimental results in Table 1.

As shown in this table, both of our LAP and GAP model are able to utilize the unlabeled data to increase the overall performance comparing with only using labeled data. Our LAP model performs slightly worse than the NMP model, which we attribute to the increased model complexity by incorporating extra encoder and decoders to deal with the latent variable. However, our LAP model achieved comparable results on semi-supervised parsing as the DPPP model, while our LAP model is simple and straightforward without additional inference procedure. Instead, the DPPP model has to sample from the posterior of the structure by using a “GUMBEL-MAX trick” to approximate the categorical distribution at each step, which is intensively computationally expensive. Further, our GAP model achieved the best results among all these methods, by successfully leveraging the the unlabeled data in an appropriate manner. We owe this success to such a fact: GAP is able to calculate the exact expectation of the arc-decomposed latent variable, *the dependency tree structure*, in the ELBO for the complete data likelihood when the data is unlabeled, rather than using sampling

Model	UAS
DPPP(Corro & Titov, 2018)(L)	88.79
DPPP(Corro & Titov, 2018)(L+U)	89.50
CRFAE(Cai et al., 2017)(L+U)	82.34
NMP(Kiperwasser & Goldberg, 2016)(L)	89.64
NTP (L)	89.63
Self-training (L+U)	87.81
LAP (L)	89.37
LAP (L+U)	<b>89.49</b>
GAP (L)	89.65
GAP (L+U)	<b>89.96</b>

Table 1: Comparing model performance on WSJ data set with 10% labeled data. “L” means only 10% labeled data is used, while “L+U” means both 10% labeled and 90% unlabeled data are used.

Model	Dutch	Spanish	English	French	Croatian	German	Italian	Russian	Japanese
NMP (L)	76.11	82.00	75.51	83.07	77.44	74.07	82.85	75.18	93.46
NTP (L)	76.20	82.09	75.57	83.12	77.51	74.13	82.99	75.23	93.54
LAP (L)	76.15	81.93	75.36	83.09	77.45	74.14	83.07	74.84	93.38
GAP (L)	76.23	81.97	75.75	83.11	77.49	74.16	83.14	75.17	93.52
CRFAE (L+U)	71.32	74.67	68.52	77.35	69.89	68.44	76.37	68.64	87.26
ST (L+U)	75.37	80.86	72.76	81.38	76.10	73.45	82.74	72.57	91.43
LAP (L+U)	76.29	82.48	75.48	83.23	77.78	74.48	83.34	75.22	93.65
GAP (L+U)	76.54	82.56	76.21	83.26	77.83	74.63	83.54	75.69	93.92

Table 2: In this table we compare different models on multiple languages from UD. Models were trained in a fully supervised fashion with labeled data only (noted as “L”) or semi-supervised (notes as “L+U”). “ST” stands for self-training.

to approximate the true expectation. Self-training using NMP with both labeled and unlabeled data is also included as a base-line, where the performance is deteriorated without appropriately using the unlabeled data.

### 7.3 SEMI-SUPERVISED DEPENDENCY PARSING ON THE UD DATA SET

We also evaluated our models on multiple languages from the UD data and compared the model performance with the semi-supervised version of CRFAE and the fully supervised NMP and NTP. To fully simulate the low-resource scenario, no external word embeddings were used.

We summarize the results in Table 2. First, when using labeled data only, LAP and GAP have similar performance as NMP and NTP. Second, we note that our LAP and GAP models do benefit from the unlabeled data, compared to using labeled data only. Both our LAP and GAP model are able to exploit the hidden information in the unlabeled data to improve the performance. Comparing between LAP and GAP, we notice GAP in general has better performance than LAP, and can better leverage the information in the unlabeled data to boost the performance. These results validate that GAP is especially useful for low-resource languages with few annotations. We also experimented using self-training on the labeled and unlabeled data with the NMP model. As results show, self-training deteriorate the performance especially when the size of the training data is small.

## 8 CONCLUSION

In this paper, we present two semi-supervised parsers, which are locally autoencoding parser (LAP) and globally autoencoding parser (GAP). Both of them are end-to-end learning systems enhanced with neural architecture, capable of utilizing the latent information within the unlabeled data together with labeled data to improve the parsing performance, without using external resources. More importantly, our GAP model outperforms the previous published (Corro & Titov, 2018) semi-supervised parsing system on the WSJ data set. We attribute this success to two reasons: First, our GAP model consists both a discriminative component and a generative component. These two components are constraining and supplementing each other such that final parsing choices are made in a checked-and-balanced manner to avoid over-fitting. Second, instead of sampling from posterior of the latent variable (the dependency tree) (Corro & Titov, 2018), our model analytically computes the expectation and marginalization of the latent variable, such that the global optima can be found for the decoder, which leads to an improved performance.



## REFERENCES

- Daniel Andor, Chris Alberti, David Weiss, Aliaksei Severyn, Alessandro Presta, Kuzman Ganchev, Slav Petrov, and Michael Collins. Globally normalized transition-based neural networks. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2016.
- Joost Bastings, Ivan Titov, Wilker Aziz, Diego Marcheggiani, and Khalil Simaan. Graph Convolutional Encoders for Syntax-aware Neural Machine Translation. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2017.
- Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating Sentences from a Continuous Space. In *Proc. International Conference on Learning Representation (ICLR)*, 2016.
- Jiong Cai, Yong Jiang, and Kewei Tu. CRF Autoencoder for Unsupervised Dependency Parsing. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2017.
- Danqi Chen and Christopher D Manning. A fast and accurate dependency parser using neural networks. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2014.
- Mingda Chen, Qingming Tang, Karen Livescu, and Kevin Gimpel. Variational sequential labelers for semi-supervised learning. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2018.
- Caio Corro and Ivan Titov. Differentiable Perturb-and-Parse: Semi-Supervised Parsing with a Structured Variational Autoencoder. In *Proc. International Conference on Learning Representation (ICLR)*, 2018.
- Hang Cui, Renxu Sun, Keya Li, Min-Yen Kan, and Tat-Seng Chua. Question Answering Passage Retrieval Using Dependency Relations. In *Proc. of International Conference on Research and Development in Information Retrieval (SIGIR)*, 2005.
- Aron Culotta and Jeffrey Sorensen. Dependency tree kernels for relation extraction. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2004.
- Marie-Catherine De Marneffe and Christopher D Manning. The Stanford typed dependencies representation. In *Proc. the International Conference on Computational Linguistics (COLING)*, 2008.
- Yuan Ding and Martha Palmer. Machine translation using probabilistic synchronous dependency insertion grammars. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2007.
- Timothy Dozat and Christopher D. Manning. Deep biaffine attention for neural dependency parsing. In *Proc. International Conference on Learning Representation (ICLR)*, April 2017.
- Timothy Dozat, Peng Qi, and Christopher D. Manning. Stanford’s graph-based neural dependency parser at the conll 2017 shared task. In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, 2017.
- Chris Dyer, Miguel Ballesteros, Wang Ling, Austin Matthews, and Noah A. Smith. Transition-Based Dependency Parsing with Stack Long Short-Term Memory. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, may 2015.
- Jason Eisner. Three New Probabilistic Models for Dependency Parsing: An Exploration. In *Proc. the International Conference on Computational Linguistics (COLING)*, 1996.
- Yong Jiang, Wenjuan Han, and Kewei Tu. Unsupervised Neural Dependency Parsing. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2016.
- Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized variational autoencoders. In *Proc. of the International Conference on Machine Learning (ICML)*, 2018.

- Yoon Kim, Chris Dyer, and Alexander Rush. Compound probabilistic context-free grammars for grammar induction. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2019a.
- Yoon Kim, Alexander Rush, Lei Yu, Adhiguna Kuncoro, Chris Dyer, and Gábor Melis. Unsupervised recurrent neural network grammars. In *Proc. of the Annual Meeting of the North American Association of Computational Linguistics (NAACL)*, 2019b.
- Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *ArXiv*, dec 2014.
- Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. In *Proc. International Conference on Learning Representation (ICLR)*, 2014.
- Eliyahu Kiperwasser and Yoav Goldberg. Semi-supervised dependency parsing using bilexical contextual features from auto-parsed data. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2015.
- Eliyahu Kiperwasser and Yoav Goldberg. Simple and accurate dependency parsing using bidirectional lstm feature representations. *Transactions of the Association for Computational Linguistics (TACL)*, 2016.
- Terry Koo, Xavier Carreras, and Michael Collins. Simple Semi-supervised Dependency Parsing. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2008.
- Sandra Kubler, Ryan McDonald, Joakim Nivre, and Graeme Hirst. *Dependency Parsing*. Morgan and Claypool Publishers, 2009.
- Zhenghua Li, Min Zhang, and Wenliang Chen. Ambiguity-aware ensemble training for semi-supervised dependency parsing. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2014.
- Yang Liu, Furu Wei, Sujian Li, Heng Ji, Ming Zhou, and Houfeng Wang. A Dependency-Based Neural Network for Relation Classification. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2015.
- Xuezhe Ma, Zecong Hu, Jingzhou Liu, Nanyun Peng, Graham Neubig, and Eduard Hovy. Stack-Pointer Networks for Dependency Parsing. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*. Association for Computational Linguistics, 2018.
- Diego Marcheggiani and Ivan Titov. Encoding Sentences with Graph Convolutional Networks for Semantic Role Labeling. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2017.
- Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of English: the penn treebank. *Computational Linguistics*, 1993.
- Ryan McDonald. *Discriminative learning and spanning tree algorithms for dependency parsing*. PhD thesis, University of Pennsylvania, 2006.
- Ryan McDonald, Koby Crammer, and Fernando Pereira. Online large-margin training of dependency parsers. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2005.
- Ryan Mcdonald, Joakim Nivre, Yvonne Quirnbach-brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Tckstrm, Claudia Bedini, Nria Bertomeu, and Castell Jungmee Lee. Universal dependency annotation for multilingual parsing. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2013.
- Yishu Miao and Phil Blunsom. Language as a latent variable: Discrete generative models for sentence compression. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2016.
- Joakim Nivre. Incrementality in deterministic dependency parsing. In *Proceedings of the Workshop on Incremental Parsing: Bringing Engineering and Cognition Together*, 2004.

- Joakim Nivre. Algorithms for deterministic incremental dependency parsing. *Computational Linguistics*, 2008.
- Joakim Nivre. The inside-outside recursive neural network model for dependency parsing. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2014.
- Mark A. Paskin. Cubic-time parsing and learning algorithms for grammatical bigram models. Technical report, EECS Department, University of California, Berkeley, 2001.
- Wenzhe Pei, Tao Ge, and Baobao Chang. An effective neural network model for graph-based dependency parsing. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2015.
- Siva Reddy, Oscar Tackstrom, Michael Collins, Tom Kwiatkowski, Dipanjan Das, Mark Steedman, and Mirella Lapata. Transforming Dependency Structures to Logical Forms for Semantic Parsing. *Transactions of the Association for Computational Linguistics (TACL)*, pp. 127–140, 2016.
- Anders Søgaard and Christian Rishøj. Semi-supervised dependency parsing using generalized tri-training. In *Proc. the International Conference on Computational Linguistics (COLING)*, 2010.
- Kihyuk Sohn, Xinchun Yan, and Honglak Lee. Learning structured output representation using deep conditional generative models. In *The Conference on Advances in Neural Information Processing Systems (NIPS)*, 2015.
- Karl Stratos. Mutual Information Maximization for Simple and Accurate Part-Of-Speech Induction. In *Proc. of the Annual Meeting of the North American Association of Computational Linguistics (NAACL)*, 2019.
- Toshiyuki Tanaka. A Theory of Mean Field Approximation. In *The Conference on Advances in Neural Information Processing Systems (NIPS)*, 1999.
- Ben Taskar, Vassil Chatalbashev, Daphne Koller, and Carlos Guestrin. Learning structured prediction models: A large margin approach. In *Proc. of the International Conference on Machine Learning (ICML)*, 2005.
- Kewei Tu and Vasant Honavar. Unambiguity regularization for unsupervised learning of probabilistic grammars. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2012.
- Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *The Conference on Advances in Neural Information Processing Systems (NIPS)*, 2017.
- David Weiss, Chris Alberti, Michael Collins, and Slav Petrov. Structured training for neural network transition-based parsing. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2015.
- Sam Wiseman and Alexander M. Rush. Sequence-to-sequence learning as beam-search optimization. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2016.
- Xiaoqing Zheng. Incremental graph-based neural dependency parsing. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2017.
- Chunting Zhou and Graham Neubig. Multi-space variational encoder-decoders for semi-supervised labeled sequence transduction. In *Proc. of the Annual Meeting of the Association Computational Linguistics (ACL)*, 2017.

## A APPENDIX

## DETAILS OF THE LAP MODEL

## ELBO OF LAP’S ORIGINAL OBJECTIVE

**Lemma A.1.**  $\mathcal{J}_{lap}$  is the ELBO (evidence lower bound) of the original objective  $\mathcal{J}$ , with an input sequence  $\mathbf{x}$ .

Denote the encoder  $Q$  is a distribution used to approximate the true posterior distribution  $P_\phi(z|\mathbf{x})$ , parameterized by  $\phi$  such that  $Q$  encoding the input into the latent space  $z$ .

*Proof.*

$$\begin{aligned} \log P_\theta(\mathbf{x})P_\omega^\epsilon(\mathcal{T}|\mathbf{x}) &= \underbrace{\log P_\theta(\mathbf{x})}_{\mathcal{U}} + \underbrace{\epsilon \log P_\omega(\mathcal{T}|\mathbf{x})}_{\mathcal{L}} \\ \mathcal{U} &= \log \int_z Q_\phi(z|\mathbf{x}) \frac{P_\theta(\mathbf{x})}{Q_\phi(z|\mathbf{x})} dz \\ &\geq \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}|z)] - \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} \left[ \log \frac{Q_\phi(z|\mathbf{x})}{P_\theta(\mathbf{x})} \right] \\ &= \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}|z)] - \mathbb{KL}(Q_\phi(z|\mathbf{x})||P_\theta(\mathbf{x})), \text{ (ELBO of traditional VAE)} \\ \mathcal{L} &= \epsilon \log P_\omega(\mathcal{T}|\mathbf{x}) \\ &= \epsilon \log \int_z P_\omega(\mathcal{T}|z) Q_\phi(z|\mathbf{x}) dz \\ &= \epsilon \log \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [P_\omega(\mathcal{T}|z)] \\ &\geq \epsilon \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\omega(\mathcal{T}|z)]. \end{aligned}$$

Combining  $\mathcal{U}$  and  $\mathcal{L}$  leads to the fact:

$$\mathcal{U} + \mathcal{L} \geq \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}|z)] - \mathbb{KL}(Q_\phi(z|\mathbf{x})||P_\theta(\mathbf{x})) + \epsilon \mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\omega(\mathcal{T}|z)] = \mathcal{J}_{lap}$$

□

In practice, similar as VAE-style models,  $\mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}|z)]$  is approximated by  $\frac{1}{N} \sum_{j=1}^N \log P_\theta(\mathbf{x}|z_j)$  and  $\mathbb{E}_{z \sim Q_\phi(z|\mathbf{x})} [\log P_\omega(\mathcal{T}|z)]$  by  $\frac{1}{N} \sum_{j=1}^N \log P_\omega(\mathcal{T}|z_j)$ , where  $z_j$  is the  $j$ th sample of  $N$  samples sampled from  $Q_\phi(z|\mathbf{x})$ . At prediction stage, we simply use  $\mu_z$  rather than sampling  $z$ .

## MEAN FIELD APPROXIMATION AND ANNEALING

Here we used a mean field approximation (Tanaka, 1999) together with the conditional independence assumption by assuming  $P_\theta(z|\mathbf{x}) \approx \prod_{t=1}^l Q_\phi(z_t|x_t)$ . The generative model  $P_\theta(\mathbf{x}|z)$  acting as decoder parameterized by  $\theta$  tries to regenerate the specific input  $x_t$  at time step  $t$  from the latent space  $z_t$ , as we assume conditional independence in the generative process among  $P_\theta(x_t|z_t)$ . The encoder and the decoder are trained jointly in the classical variational autoencoder framework, by minimizing the KL divergence between the approximated posterior and the true posterior.

We describe the encoder and decoder formulation. We parameterize the encoder  $Q_\phi(z_t|x_t)$  in such a way: First a bi-LSTM is used to obtain a non-linear transformation  $h_t$  of the original  $x_t$ ; then two separate MLPs are used to compute the mean  $\mu_{z_t}$  and the variance  $\sigma_{z_t}^2$ . The generative story  $P_\theta(x_t|z_t)$  follows such parameterization: we used a MLP of two hidden layers in-between to take  $z_t$  as the input, and then predict the word (or POS tag) over the vocabulary, such that the reconstruction probability can be measured.

Following traditional VAE training paradigms, we also apply the “re-parameterization” trick (Kingma & Welling, 2014) to circumvent the non-differentiable sampling procedure to sample  $z_t$  from the  $Q_\phi(z_t|x_t)$ . Instead of directly sample from  $\mathcal{N}(\mu_{z_t}, \sigma_{z_t}^2)$ , we form  $z_t = \mu_{z_t} + \epsilon \odot \sigma_{z_t}$  by sampling  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . In addition, to avoid hindering learning during the initial training phases, following previous works (Chen et al., 2018; Bowman et al., 2016), we anneal the temperature on the KL divergence term from a small value to 1.

## EMPIRICAL BAYESIAN TREATMENT

From an empirical Bayesian perspective, rather than fixing the prior using some certain distributions, it is beneficial to estimate the prior distribution directly from the data by treating prior’s parameters part of the model parameters. Similar to the approach used in the previous study (Chen et al., 2018), LAP also learns the priors from the data by updating them iteratively. We initialize the priors from a standard Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , where  $\mathbf{I}$  is an identity matrix. During the training, the current priors are updated using the last optimized posterior, following the rule:

$$\pi_{\theta}^k(\mathbf{z}) = \sum_{\mathbf{x}} Q_\phi^{k-1}(\mathbf{z}|\mathbf{x})P(\mathbf{x}),$$

where  $P(\mathbf{x})$  represents the empirical data distribution, and  $k$  the iteration step. Empirical Bayesian is also named as “maximum marginal likelihood”, such that our approach here is to marginalize over the missing observation as a random variable.

## INCORPORATING POS AND EXTERNAL EMBEDDINGS

In previous studies (Chen & Manning, 2014; Dozat & Manning, 2017; Dozat et al., 2017; Kiperwasser & Goldberg, 2016) exploring parsing using neural architectures, POS tags and external embeddings have been shown to contain important information characterizing the dependency relationship between a head and a child. Therefore, in addition to the variational autoencoding framework taking as input the randomly initialized word embeddings, optionally we can build the same structure for POS to reconstruct tags and for external embeddings to reconstruct words as well, whose variational objectives are  $\mathcal{U}_p$  and  $\mathcal{U}_e$  respectively. Hence, the final variational objective can be a combination of three:  $\mathcal{U} = \mathcal{U}_w$  (The original  $\mathcal{U}$  in Lemma A.1) +  $\mathcal{U}_p$  +  $\mathcal{U}_e$  (or just  $\mathcal{U} = \mathcal{U}_w + \mathcal{U}_p$  if external embeddings are not used).

## DETAILS OF THE GAP MODEL

### MARGINALIZATION AND EXPECTATION OF LATENT PARSE TREES

Assuming the sentence is of length  $l$ , and we have obtained a arc decomposed scoring matrix  $\mathbf{S}$  of size  $l \times l$ , and an entry  $\mathbf{S}[i, j]_{i \neq j, j \neq 0}$  stands for the arc score where  $i$ th word is the head and  $j$ th word the modifier. We first describe the **inside** algorithm to compute the marginalization of all possible projective trees in Algo.2.

We then describe the **outside** algorithm to compute the outside tables in Algo. 3. In this algorithm,  $\oplus$  stands for the *logaddexp* operation.

Finally, with the inside table  $\alpha$ , outside table  $\beta$  and the marginalization  $Z$  of all possible latent trees, we can compute the expectation of latent tree in an arc-decomposed manner. Algo. 4 describes the procedure. It results the matrix  $\mathbf{P}$  containing the expectation of all individual arcs by marginalize over all other arcs except itself.

Light modification is needed in our study to calculate the expectation w.r.t. the posterior distribution  $Q(\mathcal{T}) = P_{\Theta, \Phi}(\mathcal{T}|\mathbf{m}, \mathbf{x})$ , as we have

**Algorithm 2** Inside Algorithm

---

**Input:**  $S$   
**Output:**  $\alpha, Z$

```

1:  $\alpha \leftarrow -\infty$ 
2: for  $s \in 0 \dots l - 1$  do
3:   if  $s > 0$  then
4:      $\alpha[s, s, L, C] \leftarrow 0$ 
5:   end if
6:    $\alpha[s, s, R, C] \leftarrow 0$ 
7: end for
8: for  $k \in 1 \dots l - 1$  do
9:   for  $s \in 0 \dots l - k$  do
10:     $t = s + k$ 
11:    if  $s > 0$  then
12:       $\alpha[s, t, L, I] \leftarrow \log \sum_{u \in [s, t]} \exp(\alpha[s, u, L, C] + \alpha[u, t, L, C]) + \mathbf{S}[t, s]$ 
13:    end if
14:     $\alpha[s, t, R, I] \leftarrow \log \sum_{u \in [s, t]} \exp(\alpha[s, u, L, C] + \alpha[u, t, L, C]) + \mathbf{S}[s, t]$ 
15:    if  $s > 0$  then
16:       $\alpha[s, t, L, C] \leftarrow \log \sum_{u \in [s, t-1]} \exp(\alpha[s, u, L, C] + \alpha[u, t, L, C])$ 
17:    end if
18:     $\alpha[s, t, R, C] \leftarrow \log \sum_{u \in [s+1, t]} \exp(\alpha[s, u, R, I] + \alpha[u, t, R, I])$ 
19:  end for
20: end for
21:  $Z \leftarrow \alpha[0, l - 1, R, C]$ 

```

---

**Algorithm 3** Outside Algorithm

---

**Input:**  $S, \alpha$   
**Output:**  $\beta$

```

1:  $\beta \leftarrow -\infty$ 
2:  $\beta[0, l - 1, R, C] \leftarrow 0$ 
3: for  $k \in l - 1 \dots 1$  do
4:   for  $s \in 0 \dots l - k$  do
5:      $t = s + k$ 
6:      $\beta[s, s + 1 : t + 1, R, I] \leftarrow \bigoplus(\beta[s, s + 1 : t + 1, R, I], \beta[s, t, R, C] + \alpha[s + 1 : t + 1, t, R, C])$ 
7:      $\beta[s + 1 : t + 1, t, R, C] \leftarrow \bigoplus(\beta[s + 1 : t + 1, t, R, C], \beta[s, t, R, C] + \alpha[s, s + 1 : t + 1, R, I])$ 
8:     if  $s > 0$  then
9:        $\beta[s, s : t, L, C] \leftarrow \bigoplus(\beta[s, s : t, L, C], \beta[s, t, R, C] + \alpha[s : t, t, L, I])$ 
10:       $\beta[s : t, t, L, I] \leftarrow \bigoplus(\beta[s : t, t, L, I], \beta[s, t, L, C] + \alpha[s, s : t, L, C])$ 
11:     end if
12:      $\beta[s, s : t, R, C] \leftarrow \bigoplus(\beta[s, s : t, R, C], \beta[s, t, R, I] + \alpha[s + 1 : t + 1, t, L, C] + \mathbf{S}[s, t])$ 
13:      $\beta[s + 1 : t + 1, t, L, C] \leftarrow \bigoplus(\beta[s + 1 : t + 1, t, L, C], \beta[s, t, R, I] + \alpha[s, s : t, R, C] + \mathbf{S}[s, t])$ 
14:     if  $s > 0$  then
15:        $\beta[s, s : t, R, C] \leftarrow \bigoplus(\beta[s, s : t, R, C], \beta[s, t, L, I] + \alpha[s + 1 : t + 1, t, L, C] + \mathbf{S}[t, s])$ 
16:        $\beta[s + 1 : t + 1, t, L, C] \leftarrow \bigoplus(\beta[s + 1 : t + 1, t, L, C], \beta[s, t, L, I] + \alpha[s, s : t, R, C] + \mathbf{S}[t, s])$ 
17:     end if
18:   end for
19: end for

```

---

**Algorithm 4** Arc Decomposed Expectation

---

**Input:**  $\alpha, \beta, Z$   
**Output:**  $P$

- 1:  $P \leftarrow 0$
- 2: **for**  $s \in 0 \dots l - 2$  **do**
- 3:     **for**  $t \in s + 1 \dots l - 1$  **do**
- 4:         **if**  $s \neq t$  **then**
- 5:              $P[s, t] \leftarrow \exp(\alpha[s, t, R, I] + \beta[s, t, R, I] - Z)$
- 6:             **if**  $s > 0$  **then**
- 7:                  $P[t, s] \leftarrow \exp(\alpha[s, t, L, I] + \beta[s, t, L, I] - Z)$
- 8:             **end if**
- 9:         **end if**
- 10:     **end for**
- 11: **end for**

---

$$\begin{aligned}
P_{\Theta, \Phi}(\mathcal{T} | \mathbf{m}, \mathbf{x}) &= \frac{P_{\Theta, \Phi}(\mathcal{T}, \mathbf{m} | \mathbf{x})}{P_{\Theta, \Phi}(\mathbf{m} | \mathbf{x})} \\
&= \frac{\exp \sum_{(h, m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m)}{Z(\mathbf{x})} / \sum_{\mathcal{T} \in \mathbb{T}} \left[ \frac{\exp \sum_{(h, m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m)}{Z(\mathbf{x})} \right] \\
&= \frac{\exp \sum_{(h, m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m)}{Z'(\mathbf{x})},
\end{aligned}$$

where  $Z'(\mathbf{x}) = \sum_{\mathcal{T} \in \mathbb{T}} \exp \sum_{(h, m) \in \mathcal{T}} s'_{\Phi, \Theta}(h, m)$  is the real marginal we need to calculate using the transformed scoring matrix  $\mathbf{S}'$  as input in the inside algorithm. Each entry in this transformed scoring matrix is defined in the text as  $s'_{\Phi, \Theta}(h, m)$ .

CONVEXITY OF ELBO W.R.T.  $\Theta$ 

In this section we derive the strict convexity of ELBO w.r.t.  $\Theta$ . Since we only care about the term containing  $\Theta$ , the KL divergence term degenerates to a constant. For sentence  $i$ ,  $Q(\mathcal{T}_i)$  has been derived in the previous section as matrix  $\mathbf{P}$  and  $\mathbb{1}$  is the indication function.

$$\begin{aligned}
&\max_{\Theta} \sum_i \mathbb{E}_{\mathcal{T}_i \sim Q(\mathcal{T}_i)} [\log P_{\Theta}(\mathbf{m}_i | \mathcal{T}_i)] - \mathbb{KL} [Q(\mathcal{T}_i) || P_{\Phi}(\mathcal{T}_i | \mathbf{x}_i)] \\
&\max_{\Theta} \sum_i \sum_{\mathcal{T}_i \in \mathbb{T}(\mathbf{x}_i)} Q(\mathcal{T}_i) \log P_{\Theta}(\mathbf{m}_i | \mathcal{T}_i) + Const \\
&\max_{\Theta} \sum_{(h \rightarrow m)} \log \theta_{mh} \mathbb{E}_{(h \rightarrow m) \sim Q} \mathbb{1}(h \rightarrow m) \\
&\max_{\Theta} \sum_{(h \rightarrow m)} Q(\mathbb{1}(h \rightarrow m)) \log \theta_{mh} \quad Q(\mathbb{1}(h \rightarrow m)) \text{ is a Bernoulli distribution, indicating whether the arc } (h \rightarrow m) \text{ exists.} \\
&s.t. \sum_m \theta_{mh} = 1 \quad \forall h.
\end{aligned} \tag{3}$$

## DATA SET STATISTICS

We show the details of the statistics of the WSJ data set, which is the Stanford Dependency conversion (De Marneffe & Manning, 2008) of the Penn Treebank (Marcus et al., 1993) and the statistics of the languages we used in UD (Universal Dependency) 2.3 (McDonald et al., 2013) here.

Language	WSJ	Dutch	Spanish	English	French	Croatian	German	Italian	Russian	Japanese
Training	39832	12269	14187	2914	14450	6983	13814	13121	3850	7133
Development	1700	718	1400	707	1476	849	799	564	579	511
Testing	2416	596	426	769	416	1057	977	482	601	551

Table 3: Statistics of multiple languages we used in our experiments are shown here. The table shows number of sentences in the training, development and test data divisions.