# Globally Normalized Transition-Based Neural Networks

**Daniel Andor, Chris Alberti, David Weiss, Aliaksei Severyn, Alessandro Presta,**
**Kuzman Ganchev, Slav Petrov and Michael Collins**
Google Inc
New York, NY
{andor,chrisalberti,djweiss,severyn,apresta,kuzman,slav,mjcollins}@google.com

## Abstract

We introduce a globally normalized transition-based neural network model that achieves state-of-the-art part-of-speech tagging, dependency parsing and sentence compression results. Our model is a simple feed-forward neural network that operates on a task-specific transition system, yet achieves comparable or better accuracies than recurrent models. The key insight is based on a novel proof illustrating the label bias problem and showing that globally normalized models can be strictly more expressive than locally normalized models.

## 1 Introduction

Neural network approaches have taken the field of natural language processing (NLP) by storm. In particular, variants of long short-term memory (LSTM) networks (Hochreiter and Schmidhuber, 1997) have produced impressive results on some of the classic NLP tasks such as part-of-speech tagging (Ling et al., 2015), syntactic parsing (Vinyals et al., 2015) and semantic role labeling (Zhou and Xu, 2015). One might speculate that it is the recurrent nature of these models that enables these results.

In this work we demonstrate that simple feed-forward networks without any recurrence can achieve comparable or better accuracies than LSTMs, as long as they are globally normalized. Our model, described in detail in Section 2, uses a transition system (Nivre, 2006) and feature embeddings as introduced by Chen and Manning (2014). We do not use any recurrence, but perform beam search for maintaining multiple hypotheses and intro-

duce global normalization with a conditional random field (CRF) objective (Bottou et al., 1997; Le Cun et al., 1998; Lafferty et al., 2001) to overcome the label bias problem that locally normalized models suffer from. Since we use beam inference, we approximate the partition function by summing over the elements in the beam, and use early updates (Collins and Roark, 2004; Zhou et al., 2015). We compute gradients based on this approximate global normalization and perform full backpropagation training of all neural network parameters based on the CRF loss.

We revisit the label bias problem in Section 3 and provide a novel proof that globally normalized models are strictly more expressive than locally normalized models. Lookahead features can partially mitigate this discrepancy, but cannot fully compensate for it—a point to which we return later. To empirically demonstrate the effectiveness of global normalization, we evaluate our model on part-of-speech tagging, syntactic dependency parsing and sentence compression (Section 4). Our model achieves state-of-the-art accuracy on all of these tasks, matching or outperforming LSTMs while being significantly faster. In particular for dependency parsing on the Wall Street Journal we achieve the best-ever published unlabeled attachment score of 94.41%.

As discussed in more detail in Section 5, we also outperform previous structured training approaches used for neural network transition-based parsing. Our ablation experiments show that we outperform Weiss et al. (2015) and Alberti et al. (2015) because we do global backpropagation training of all model parameters, while they fix the neural network parameters when training the global part of their model. We also outperform Zhou et al. (2015) despite using a smaller beam. To shed additional light on the label bias problem in practice, we provide a sentence

compression example where the local model completely fails. We then demonstrate that a globally normalized parsing model without any lookahead features is almost as accurate as our best model, while a locally normalized model loses more than 10% absolute in accuracy because it cannot effectively incorporate evidence as it becomes available.

## 2 Model

At its core, our model is an incremental transition-based parser (Nivre, 2006). To apply it to different tasks we only need to adjust the transition system and the input features.

### 2.1 Transition System

Given an input $x$, most often a sentence, we define:
- A set of states $\mathcal{S}$.
- A special start state $s^\dagger \in \mathcal{S}$.
- A set of allowed decisions $\mathcal{A}(s)$ for all $s \in \mathcal{S}$.
- A transition function $t(s, d)$ returning a new state $s'$ for any decision $d \in \mathcal{A}(s)$.

We drop the dependence on $x$ for brevity. We will use a function $\rho(s, d; \theta)$ to compute the score of decision $d$ in state $s$. The vector $\theta$ contains the model parameters and we assume that $\rho(s, d; \theta)$ is differentiable with respect to $\theta$.

Throughout this work we will use transition systems in which all complete structures for the same input $x$ have the same number of decisions $n(x)$ (or $n$ for brevity). In dependency parsing for example, this is true for both the *arc-standard* and *arc-eager* transition systems (Nivre, 2006), where for a sentence $x$ of length $m$, the number of decisions for any complete parse is $n(x) = 2 \times m$.[1] A complete structure is then a sequence of decision/state pairs $(s_1, d_1) \ldots (s_n, d_n)$ such that $s_1 = s^\dagger$, $d_i \in \mathcal{S}(s_i)$ for $i = 1 \ldots n$, and $s_{i+1} = t(s_i, d_i)$. We use the notation $d_{1:j}$ to refer to a decision sequence $d_1 \ldots d_j$.

We assume that there is a one-to-one mapping between decision sequences $d_{1:j}$ and states $s_j$: that is, we essentially assume that a state encodes the entire history of decisions. Thus, each state can be reached by a unique decision sequence from $s^\dagger$.[2] We will use decision sequences $d_{1:j}$ and states interchangeably: in a slight abuse of notation, we

define $\rho(d_{1:j}, d; \theta)$ to be equal to $\rho(s, d; \theta)$ where $s$ is the state reached by decisions $d_{1:j}$.

The scoring function $\rho(s, d; \theta)$ can be defined in a number of ways. In this work, following Chen and Manning (2014), Weiss et al. (2015), and Zhou et al. (2015), we define it via a feed-forward neural network as

$$\rho(s, d; \theta) = \phi(s; \theta^{(l)}) \cdot \theta^{(d)}.$$

Here $\theta^{(l)}$ are the parameters of the neural network, excluding the parameters at the final layer. $\theta^{(d)}$ are the final layer parameters for decision $d$. $\phi(s; \theta^{(l)})$ is the representation for state $s$ computed by the neural network under parameters $\theta^{(l)}$. Note that the score is linear in the parameters $\theta^{(d)}$. We next describe how softmax-style normalization can be performed at the local or global level.

### 2.2 Global vs. Local Normalization

In the Chen and Manning (2014) style of greedy neural network parsing, the conditional probability distribution over decisions $d_j$ given context $d_{1:j-1}$ is defined as

$$p(d_j | d_{1:j-1}; \theta) = \frac{\exp \rho(d_{1:j-1}, d_j; \theta)}{Z_L(d_{1:j-1}; \theta)}, \quad (1)$$

where

$$Z_L(d_{1:j-1}; \theta) = \sum_{d' \in \mathcal{A}(d_{1:j-1})} \exp \rho(d_{1:j-1}, d'; \theta).$$

Each $Z_L(d_{1:j-1}; \theta)$ is a *local* normalization term. The probability of a sequence of decisions $d_{1:n}$ is

$$p_L(d_{1:n}) = \prod_{j=1}^{n} p(d_j | d_{1:j-1}; \theta)$$
$$= \frac{\exp \sum_{j=1}^{n} \rho(d_{1:j-1}, d_j; \theta)}{\prod_{j=1}^{n} Z_L(d_{1:j-1}; \theta)}. \quad (2)$$

Beam search can be used to attempt to find the maximum of (2) with respect to $d_{1:n}$.

In contrast, a Conditional Random Field (CRF) defines a distribution $p_G(d_{1:n})$ as follows:

$$p_G(d_{1:n}) = \frac{\exp \sum_{j=1}^{n} \rho(d_{1:j-1}, d_j; \theta)}{Z_G(\theta)}, \quad (3)$$

where

$$Z_G(\theta) = \sum_{d'_{1:n} \in \mathcal{D}_n} \exp \sum_{j=1}^{n} \rho(d'_{1:j-1}, d'_j; \theta)$$

---

[1]Note that this is not true for the *swap* transition system defined in Nivre (2009).

[2]It is straightforward to extend the approach to make use of dynamic programming in the case where the same state can be reached by multiple decision sequences.

and $\mathcal{D}_n$ is the set of all valid sequences of decisions of length $n$. $Z_G(\theta)$ is a *global* normalization term. The inference problem is now to find

$$\operatorname*{argmax}_{d_{1:n} \in \mathcal{D}_n} p_G(d_{1:n}) = \operatorname*{argmax}_{d_{1:n} \in \mathcal{D}_n} \sum_{j=1}^{n} \rho(d_{1:j-1}, d_j; \theta).$$

Beam search can again be used to approximately find the argmax.

## 2.3 Training

Training data consists of inputs $x$ paired with gold decision sequences $d_{1:n}^*$. We use stochastic gradient descent on the negative log-likelihood of the data under the model. Under a locally normalized model, the negative log-likelihood is

$$L_{\text{local}}(d_{1:n}^*; \theta) = -\ln p_L(d_{1:n}^*; \theta) = \qquad (4)$$
$$-\sum_{j=1}^{n} \rho(d_{1:j-1}^*, d_j^*; \theta) + \sum_{j=1}^{n} \ln Z_L(d_{1:j-1}^*; \theta),$$

whereas under a globally normalized model it is

$$L_{\text{global}}(d_{1:n}^*; \theta) = -\ln p_G(d_{1:n}^*; \theta) =$$
$$-\sum_{j=1}^{n} \rho(d_{1:j-1}^*, d_j^*; \theta) + \ln Z_G(\theta). \qquad (5)$$

A significant practical advantange of the locally normalized cost (4) is that it factorizes into $n$ independent terms, each of which can be computed exactly and minimized separately. By contrast, the $Z_G$ term in (5) contains a sum over $d_{1:n}' \in \mathcal{D}_n$ that is in many cases intractable.

To make learning tractable with the globally normalized model, we use beam search and early updates (Collins and Roark, 2004; Zhou et al., 2015). As the training sequence is being decoded, we keep track of the location of the gold path in the beam. If the gold path falls out of the beam at step $j$, a stochastic gradient step is taken on the following objective:

$$L_{\text{global-beam}}(d_{1:j}^*; \theta) =$$
$$-\sum_{i=1}^{j} \rho(d_{1:i-1}^*, d_i^*; \theta) + \ln \sum_{d_{1:j}' \in \mathcal{B}_j} \exp \sum_{i=1}^{j} \rho(d_{1:i-1}', d_i'; \theta). (6)$$

Here the set $\mathcal{B}_j$ contains all paths in the beam at step $j$, together with the gold path prefix $d_{1:j}^*$. It is straightforward to derive gradients of the loss in (6) and to back-propagate gradients to all levels of a neural network defining the score $\rho(s, d; \theta)$. If the gold path remains in the beam throughout decoding, a gradient step is performed using $\mathcal{B}_n$, the beam at the end of decoding.

## 3 The Label Bias Problem

Intuitively, we would like the model to be able to revise an earlier decision made during search, when later evidence becomes available that rules out the earlier decision as incorrect. At first glance, it might appear that a locally normalized model used in conjunction with beam search or exact search is able to revise earlier decisions. However the label bias problem (see Lafferty et al. (2001), Bottou (1991), Bottou and LeCun (2005)) means that locally normalized models often have a very weak ability to revise earlier decisions.

This section gives a more formal perspective on the label bias problem than in previous work, through a proof that globally normalized models are strictly more expressive than locally normalized models. The proof makes use of an example that gives an illustration of the label bias problem.

**Global Models can be Strictly More Expressive than Local Models** Consider a tagging problem where the task is to map an input sequence $x_{1:n}$ to a decision sequence $d_{1:n}$. First, consider a locally normalized model where we restrict the scoring function to access only the first $i$ input symbols $x_{1:i}$ when scoring decision $d_i$. We will return to this restriction soon. The scoring function $\rho$ can be an otherwise arbitrary function of the tuple $\langle d_{1:i-1}, d_i, x_{1:i} \rangle$:

$$p_L(d_{1:n}|x_{1:n}) = \prod_{i=1}^{n} p_L(d_i|d_{1:i-1}, x_{1:i})$$
$$= \frac{\exp \sum_{i=1}^{n} \rho(d_{1:i-1}, d_i, x_{1:i})}{\prod_{i=1}^{n} Z_L(d_{1:i-1}, x_{1:i})}.$$

Second, consider a globally normalized model

$$p_G(d_{1:n}|x_{1:n}) = \frac{\exp \sum_{i=1}^{n} \rho(d_{1:i-1}, d_i, x_{1:i})}{Z_G(x_{1:n})}.$$

This model again makes use of a scoring function $\rho(d_{1:i-1}, d_i, x_{1:i})$ restricted to the first $i$ input symbols when scoring decision $d_i$.

Define $\mathcal{P}_L$ to be the set of all possible distributions $p_L(d_{1:n}|x_{1:n})$ under the local model obtained as the scores $\rho$ vary. Similarly, define $\mathcal{P}_G$ to be the set of all possible distributions $p_G(d_{1:n}|x_{1:n})$ under the global model. Here a "distribution" is a function from a pair $(x_{1:n}, d_{1:n})$ to a probability $p(d_{1:n}|x_{1:n})$. Our main result is the following:

**Theorem 3.1**
$\mathcal{P}_L$ *is a strict subset of* $\mathcal{P}_G$*, that is* $\mathcal{P}_L \subsetneq \mathcal{P}_G$.

To prove this we will first prove that $\mathcal{P}_L \subseteq \mathcal{P}_G$. This step is straightforward. We then show that $\mathcal{P}_G \nsubseteq \mathcal{P}_L$; that is, there are distributions in $\mathcal{P}_G$ that are not in $\mathcal{P}_L$. The proof that $\mathcal{P}_G \nsubseteq \mathcal{P}_L$ gives a clear illustration of the label bias problem.

*Proof that $\mathcal{P}_L \subseteq \mathcal{P}_G$:* We need to show that for any locally normalized distribution $p_L$, we can construct a globally normalized model $p_G$ such that $p_G = p_L$. Consider a locally normalized model with scores $\rho(d_{1:i-1}, d_i, x_{1:i})$. Define a global model $p_G$ with scores

$$\rho'(d_{1:i-1}, d_i, x_{1:i}) = \log p_L(d_i | d_{1:i-1}, x_{1:i}).$$

Then it is easily verified that

$$p_G(d_{1:n} | x_{1:n}) = p_L(d_{1:n} | x_{1:n})$$

for all $x_{1:n}, d_{1:n}$. □

In proving $\mathcal{P}_G \nsubseteq \mathcal{P}_L$ we will use a simple problem where every example seen in training or test data is one of the following two tagged sentences:

$$x_1 x_2 x_3 = \text{a b c}, \quad d_1 d_2 d_3 = \text{A B C}$$
$$x_1 x_2 x_3 = \text{a b e}, \quad d_1 d_2 d_3 = \text{A D E} \quad (7)$$

Note that the input $x_2 = \text{b}$ is ambiguous: it can take tags B or D. This ambiguity is resolved when the next input symbol, c or e, is observed.

Now consider a globally normalized model, where the scores $\rho(d_{1:i-1}, d_i, x_{1:i})$ are defined as follows. Define $\mathcal{T}$ as the set $\{(A, B), (B, C), (A, D), (D, E)\}$ of bigram tag transitions seen in the data. Similarly, define $\mathcal{E}$ as the set $\{(a, A), (b, B), (c, C), (b, D), (e, E)\}$ of (word, tag) pairs seen in the data. We define

$$\rho(d_{1:i-1}, d_i, x_{1:i}) \quad (8)$$
$$= \alpha \times [\![(d_{i-1}, d_i) \in \mathcal{T}]\!] + \alpha \times [\![(x_i, d_i) \in \mathcal{E}]\!]$$

where $\alpha$ is the single scalar parameter of the model, and $[\![\pi]\!] = 1$ if $\pi$ is true, 0 otherwise.

*Proof that $\mathcal{P}_G \nsubseteq \mathcal{P}_L$:* We will construct a globally normalized model $p_G$ such that there is no locally normalized model such that $p_L = p_G$.

Under the definition in (8), it is straightforward to show that

$$\lim_{\alpha \to \infty} p_G(\text{A B C}|\text{a b c}) = \lim_{\alpha \to \infty} p_G(\text{A D E}|\text{a b e}) = 1.$$

In contrast, under *any* definition for $\rho(d_{1:i-1}, d_i, x_{1:i})$, we must have

$$p_L(\text{A B C}|\text{a b c}) + p_L(\text{A D E}|\text{a b e}) \le 1 \quad (9)$$

This follows because $p_L(\text{A B C}|\text{a b c}) = p_L(\text{A}|\text{a}) \times p_L(\text{B}|\text{A, a b}) \times p_L(\text{C}|\text{A B, a b c})$ and $p_L(\text{A D E}|\text{a b e}) = p_L(\text{A}|\text{a}) \times p_L(\text{D}|\text{A, a b}) \times p_L(\text{E}|\text{A D, a b e})$. The inequality $p_L(\text{B}|\text{A, a b}) + p_L(\text{D}|\text{A, a b}) \le 1$ then immediately implies (9).

It follows that for sufficiently large values of $\alpha$, we have $p_G(\text{A B C}|\text{a b c}) + p_G(\text{A D E}|\text{a b e}) > 1$, and given (9) it is impossible to define a locally normalized model with $p_L(\text{A B C}|\text{a b c}) = p_G(\text{A B C}|\text{a b c})$ and $p_L(\text{A D E}|\text{a b e}) = p_G(\text{A D E}|\text{a b e})$. □

Under the restriction that scores $\rho(d_{1:i-1}, d_i, x_{1:i})$ depend only on the first $i$ input symbols, the globally normalized model is still able to model the data in (7), while the locally normalized model fails (see Eq. 9). The ambiguity at input symbol b is naturally resolved when the next symbol (c or e) is observed, but the locally normalized model is not able to revise its prediction.

It is easy to fix the locally normalized model for the example in (7) by allowing scores $\rho(d_{1:i-1}, d_i, x_{1:i+1})$ that take into account the input symbol $x_{i+1}$. Such lookahead is common in practice, but insufficient in general. For every amount of lookahead $k$, we can construct examples that cannot be modeled with a locally normalized model by duplicating the middle input b in (7) $k + 1$ times. Only a local model with scores $\rho(d_{1:i-1}, d_i, x_{1:n})$ that considers the entire input can capture any distribution $p(d_{1:n}|x_{1:n})$: in this case the decomposition $p_L(d_{1:n}|x_{1:n}) = \prod_{i=1}^{n} p_L(d_i | d_{1:i-1}, x_{1:n})$ makes no independence assumptions.

However, increasing the amount of context used as input comes at a cost, requiring more powerful learning algorithms, and potentially more training data. For a detailed analysis of the trade-offs between structural features in CRFs and more powerful local classifiers without structural constraints, see Liang et al. (2008); in these experiments local classifiers are unable to reach the performance of CRFs on problems such as parsing and named entity recognition where structural constraints are important. Note that there is nothing to preclude an approach that makes use of both global normalization and more powerful scoring functions $\rho(d_{1:i-1}, d_i, x_{1:n})$, obtaining the best of both worlds. The experiments that follow make use of both.

| Method | En WSJ | En-Union | | | CoNLL '09 | | | | | | | Avg - |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | News | Web | QTB | Ca | Ch | Cz | En | Ge | Ja | Sp | |
| Linear CRF | 97.17 | 97.60 | 94.58 | 96.04 | 98.81 | 94.45 | 98.90 | 97.50 | 97.14 | 97.90 | 98.79 | 97.17 |
| Ling et al. (2015) | **97.78** | 97.44 | 94.03 | 96.18 | 98.77 | 94.38 | 99.00 | 97.60 | **97.84** | 97.06 | 98.71 | 97.16 |
| Our Local (B=1) | 97.44 | 97.66 | 94.46 | 96.59 | 98.91 | 94.56 | 98.96 | 97.36 | 97.35 | 98.02 | 98.88 | 97.29 |
| Our Local (B=8) | 97.45 | 97.69 | 94.46 | 96.64 | 98.88 | 94.56 | 98.96 | 97.40 | 97.35 | 98.02 | 98.89 | 97.30 |
| Our Global (B=8) | 97.44 | **97.77** | **94.80** | **96.86** | **99.03** | **94.72** | **99.02** | **97.65** | 97.52 | **98.37** | **98.97** | **97.47** |

Table 1: Final POS tagging test set results on English WSJ and Treebank Union as well as CoNLL'09.

## 4 Experiments

To demonstrate the flexibility and modeling power of our approach, we provide experimental results on a diverse set of structured prediction tasks. We first direct our attention to POS tagging, then to syntactic dependency parsing and finally to sentence compression.

While directly optimizing the global model (5) works well, we found that training the model in two steps achieves the same precision much faster: we first pretrain the network using the local objective (4), and then perform additional training steps using the global objective (6). We pretrain all layers except the softmax layer in this way. We purposefully abstain from complicated hand engineering of input features, which might improve performance further (Durrett and Klein, 2015).

### 4.1 Part of Speech Tagging

Part of speech (POS) tagging is a classic NLP task, where modeling the structure of the output is important for achieving state-of-the-art performance.

**Data & Evaluation.** We conducted experiments on a number of different datasets: (1) English Wall Street Journal (WSJ) part of the Penn Treebank (Marcus et al., 1993) with standard POS tagging splits; (2) English "Treebank Union" multi-domain corpus containing data from the OntoNotes corpus version 5 (Hovy et al., 2006), the English Web Treebank (Petrov and McDonald, 2012), and the updated and corrected Question Treebank (Judge et al., 2006) with identical setup to Weiss et al. (2015); and (3) CoNLL '09 multilingual shared task (Hajič et al., 2009).

**Model Configuration.** Inspired by the integrated POS tagging and parsing transition system of Bohnet and Nivre (2012), we employ a simple transition system that uses only a SHIFT action and predicts the POS tag of the current word on the buffer as it gets shifted to the stack. We extract the following features on a window $\pm 3$ tokens centered at the current focus token: word, cluster, character n-gram up to length 3. We also extract the tag predicted for the previous 4 tokens. The network in these experiments has a single hidden layer with 256 units on WSJ and Treebank Union and 64 on CoNLL'09.

**Results.** In Table 1 we compare our model to a linear CRF and to the compositional character-to-word LSTM model of Ling et al. (2015). The CRF is a first-order linear model with exact inference and the same emission features as our model. It additionally also has transition features of the word, cluster and character n-gram up to length 3 on both endpoints of the transition. The results for Ling et al. (2015) were solicited from the authors.

Our local model already compares favorably against these methods on average. Using beam search with a locally normalized model does not help, but with global normalization it leads to a 7% reduction in relative error, empirically demonstrating the effect of label bias. It is also interesting to note that the set of character ngrams feature is very important, increasing average accuracy on the CoNLL'09 datasets by about 0.5% absolute. This shows that character-level modeling can also be done with a simple feed-forward netowork without recurrence.

### 4.2 Dependency Parsing

In dependency parsing the goal is to produce a directed tree representing the syntactic structure of the input sentence.

**Data & Evaluation.** We use the same corpora as in our POS tagging experiments, except that we use the standard parsing splits of the WSJ. We convert the English constituency trees to Stanford style dependencies (De Marneffe et al., 2006) using version 3.3.0 of the converter. For English, we use predicted POS tags (the same POS tags are used for all models) and exclude punctua-

| Method | WSJ | | Union-News | | Union-Web | | Union-QTB | |
|---|---|---|---|---|---|---|---|---|
| | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS |
| Martins et al. (2013) | 92.89 | 90.55 | 93.10 | 91.13 | 88.23 | 85.04 | 94.21 | 91.54 |
| Zhang and McDonald (2014) | 93.22 | 91.02 | 93.32 | 91.48 | 88.65 | 85.59 | 93.37 | 90.69 |
| Weiss et al. (2015) | 93.99 | 92.05 | 93.91 | 92.25 | 89.29 | 86.44 | 94.17 | 92.06 |
| Alberti et al. (2015) | 94.23 | 92.36 | 94.10 | 92.55 | 89.55 | 86.85 | 94.74 | 93.04 |
| Our Local (B=1) | 93.17 | 91.18 | 93.11 | 91.46 | 88.42 | 85.58 | 92.49 | 90.38 |
| Our Local (B=32) | 93.58 | 91.66 | 93.65 | 92.03 | 88.96 | 86.17 | 93.22 | 91.17 |
| Our Global (B=32) | **94.41** | **92.55** | **94.44** | **92.93** | **90.17** | **87.54** | **95.40** | **93.64** |

Table 2: Final English dependency parsing test set results (without tri-training for any method).

| Method | Catalan | | Chinese | | Czech | | English | | German | | Japanese | | Spanish | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS | UAS | LAS |
| Best Shared Task Result | - | 87.86 | - | 79.17 | - | 80.38 | - | 89.88 | - | 87.48 | - | 92.57 | - | 87.64 |
| Ballesteros et al. (2015) | 90.22 | 86.42 | 80.64 | 76.52 | 79.87 | 73.62 | 90.56 | 88.01 | 88.83 | 86.10 | 93.47 | 92.55 | 90.38 | 86.59 |
| Zhang and McDonald (2014) | 91.41 | 87.91 | 82.87 | 78.57 | 86.62 | 80.59 | 92.69 | 90.01 | 89.88 | 87.38 | 92.82 | 91.87 | 90.82 | 87.34 |
| Lei et al. (2014) | 91.33 | 87.22 | 81.67 | 76.71 | 88.76 | 81.77 | 92.75 | 90.00 | 90.81 | 87.81 | **94.04** | 91.84 | 91.16 | 87.38 |
| Bohnet and Nivre (2012) | 92.44 | 89.60 | 82.52 | 78.51 | 88.82 | 83.73 | 92.87 | 90.60 | **91.37** | **89.38** | 93.67 | 92.63 | 92.24 | 89.60 |
| Alberti et al. (2015) | 92.31 | 89.17 | 83.57 | 79.90 | 88.45 | 83.57 | 92.70 | 90.56 | 90.58 | 88.20 | 93.99 | **93.10** | 92.26 | 89.33 |
| Our Local (B=1) | 91.24 | 88.21 | 81.29 | 77.29 | 85.78 | 80.63 | 91.44 | 89.29 | 89.12 | 86.95 | 93.71 | 92.85 | 91.01 | 88.14 |
| Our Local (B=16) | 91.91 | 88.93 | 82.22 | 78.26 | 86.25 | 81.28 | 92.16 | 90.05 | 89.53 | 87.4 | 93.61 | 92.74 | 91.64 | 88.88 |
| Our Global (B=16) | **92.67** | **89.83** | **84.72** | **80.85** | **88.94** | **84.56** | **93.22** | **91.23** | 90.91 | 89.15 | 93.65 | 92.84 | **92.62** | **89.95** |

Table 3: Final CoNLL '09 dependency parsing test set results.

tion from the evaluation, as is standard. For the CoNLL '09 datasets we follow standard practice and include all punctuation in the evaluation. We follow Alberti et al. (2015) and use our own predicted POS tags so that we can include a k-best tag feature (see below) but use the supplied predicted morphological features. We report unlabeled and labeled attachment scores (UAS/LAS).

**Model Configuration.** Our model configuration is basically the same as the one originally proposed by Chen and Manning (2014) and then refined by Weiss et al. (2015). In particular, we use the arc-standard transition system and extract the same set of features as prior work: words, part of speech tags, and dependency arcs and labels in the surrounding context of the state, as well as k-best tags as proposed by Alberti et al. (2015). We use two hidden layers of 1,024 dimensions each.

**Results.** Tables 2 and Table 3 show our final parsing results and a comparison to the best systems from the literature. We obtain the best ever published results on almost all datasets, including the WSJ. The results in Table 2 are without tri-training. When we use tri-training, our WSJ accuracy improves to 94.61/92.78 (UAS/LAS), which compares favorably to the 94.26/92.41 reported by Weiss et al. (2015) with tri-training. As we

show in Section 5, these gains can be attributed to the full backpropagation training that differentiates our approach from that of Weiss et al. (2015) and Alberti et al. (2015). Our results also significantly outperform the LSTM-based approaches of Dyer et al. (2015) and Ballesteros et al. (2015).

### 4.3 Sentence Compression

Our final structured prediction task is extractive sentence compression.

**Data & Evaluation.** We follow Filippova et al. (2015), where a large news collection is used to heuristically generate compression instances. Our final corpus contains about 2.3M compression instances: we use 2M examples for training, 130k for development and 160k for the final test. We report per-token F1 score and per-sentence accuracy (A), i.e. percentage of instances that fully match the golden compressions. Following Filippova et al. (2015) we also run a human evaluation on 200 sentences where we ask the raters to score compressions for *readability* (`read`) and *informativeness* (`info`) on a scale from 0 to 5.

**Model Configuration.** The transition system for sentence compression is similar to POS tagging: we scan sentences from left-to-right and la-

| Method | Generated corpus | | Human eval | |
| --- | --- | --- | --- | --- |
| | A | F1 | read | info |
| Filippova et al. (2015) | **35.36** | **82.83** | 4.66 | 4.03 |
| Automatic | - | - | 4.31 | 3.77 |
| Our Local (B=1) | 30.51 | 78.72 | 4.58 | 4.03 |
| Our Local (B=8) | 31.19 | 75.69 | - | - |
| Our Global (B=8) | 35.16 | 81.41 | **4.67** | **4.07** |

Table 4: Sentence compression results on News data. *Automatic* refers to application of the same automatic extraction rules used to generate the News training corpus.

| Method | UAS | LAS |
| --- | --- | --- |
| Local (B=1) | 92.85 | 90.59 |
| Local (B=16) | 93.32 | 91.09 |
| Global (B=16) $\{\theta^{(d)}\}$ | 93.45 | 91.21 |
| Global (B=16) $\{W_2, \theta^{(d)}\}$ | 94.01 | 91.77 |
| Global (B=16) $\{W_1, W_2, \theta^{(d)}\}$ | 94.09 | 91.81 |
| Global (B=16) (full) | 94.38 | 92.17 |

Table 5: WSJ dev set scores for successively deeper levels of backpropagation. The *full* parameter set corresponds to backpropagation all the way to the embeddings. $W_i$: hidden layer $i$ weights.

bel each token as *keep* or *drop*. We extract features from words, POS tags, and dependency labels from a window of tokens centered on the input, as well as features from the history of predictions. We use a single hidden layer of size 400.

**Results.** Table 4 shows our sentence compression results. Our globally normalized model again significantly outperforms the local model. Beam search with a locally normalized model suffers from severe label bias issues that we discuss on a concrete example in Section 5. We also compare to the best sentence compression system from Filippova et al. (2015), a 3-layer stacked LSTM which uses dependency label information. The LSTM and our global model perform on par on both the automatic evaluation as well as the human ratings, but our model is roughly $100\times$ faster. All compressions kept approximately 42% of the tokens on average and all the models are significantly better than the automatic extractions ($p < 0.05$).

## 5 Discussion

We derived a proof for the label bias problem and the advantages of global models. We then emprirically verified this theoretical superiority by demonstrating state-of-the-art performance on three different tasks. Our experiments showed consistent improvements in accuracy for globally normalized models over locally normalized models with beam search. In this section we situate and compare our model to previous work and provide two examples of the label bias problem in practice.

### 5.1 Related Neural CRF Work

Neural network models have been been combined with conditional random fields and globally normalized models before. Bottou et al. (1997) and Le Cun et al. (1998) describe global training of neural network models for structured prediction problems. Peng et al. (2009) add a non-linear neural network layer to a linear-chain CRF and Do and Artires (2010) apply a similar approach to more general Markov network structures. Yao et al. (2014) and Zheng et al. (2015) introduce recurrence into the model and Huang et al. (2015) finally combine CRFs and LSTMs. These neural CRF models are limited to sequence labeling tasks where exact inference is possible, while our model works well when exact inference is intractable.

### 5.2 Related Transition-Based Parsing Work

For early work on neural-networks for transition-based parsing, see Henderson (2003; 2004). Our work is closest to the work of Weiss et al. (2015), Zhou et al. (2015) and Watanabe and Sumita (2015); in these approaches global normalization is added to the local model of Chen and Manning (2014). Empirically, Weiss et al. (2015) achieves the best performance, even though their model keeps the parameters of the locally normalized neural network fixed and only trains a perceptron that uses the activations as features. Their model is therefore limited in its ability to revise the predictions of the locally normalized model. In Table 5 we show that full backpropagation training all the way to the word embeddings is very important and significantly contributes to the performance of our model. We also compared training under the CRF objective with a Perceptron-like hinge loss between the gold and best elements of the beam. When we limited the backpropagation depth to training only the top layer $\theta^{(d)}$, we found negligible differences in accuracy: 93.20% and 93.28% for the CRF objective and hinge loss respectively. However,

| Method | Predicted compression | $p_L$ | $p_G$ |
|---|---|---|---|
| Local (B=1) | In Pakistan, former leader **Pervez Musharraf has appeared in court** for the first time, on treason charges. | 0.13 | 0.05 |
| Local (B=8) | In Pakistan, former leader **Pervez Musharraf has appeared in court** for the first time, on treason charges. | 0.16 | $< 10^{-4}$ |
| Global (B=8) | In Pakistan, former leader **Pervez Musharraf has appeared** in court for the first time, **on** treason charges. | 0.06 | 0.07 |

Table 6: Example sentence compressions where the label bias of the locally normalized model leads to a breakdown during beam search. The probability of each compression under the local ($p_L$) and global ($p_G$) models shows that only the global model can properly represent zero probability for the empty compression.

when training with full backpropagation the CRF accuracy is 0.2% higher and training converged more than 4× faster.

Zhou et al. (2015) perform full backpropagation training like us, but even with a much larger beam, their performance is significantly lower than ours. We also apply our model to two additional tasks, while they experiment only with dependency parsing. Finally, Watanabe and Sumita (2015) introduce recurrent components and additional techniques like max-violation updates for a corresponding constituency parsing model. In contrast, our model does not require any recurrence or specialized training.

### 5.3 Label Bias in Practice

We observed several instances of severe label bias in the sentence compression task. Although using beam search with the local model outperforms greedy inference on average, beam search leads the local model to occasionally produce empty compressions (Table 6). It is important to note that these are *not* search errors: the empty compression has higher probability under $p_L$ than the prediction from greedy inference. However, the more expressive globally normalized model does not suffer from this limitation, and correctly gives the empty compression almost zero probability.

We also present some empirical evidence that the label bias problem is severe in parsing. We trained models where the scoring functions in parsing at position $i$ in the sentence are limited to considering only tokens $x_{1:i}$; hence unlike the full parsing model, there is no ability to look ahead in the sentence when making a decision.[3] The result for a greedy model under this constraint is 76.96% UAS; for a locally normalized model with beam search is 81.35%; and for a globally normalized model is 93.60%. Thus the globally normalized model gets very close to the perfor-

mance of a model with full lookahead, while the locally normalized model with a beam gives dramatically lower performance. In our final experiments with full lookahead, the globally normalized model achieves 94.01% accuracy, compared to 93.07% accuracy for a local model with beam search. Thus adding lookahead allows the local model to close the gap in performance to the global model; however there is still a significant difference in accuracy, which may in large part be due to the label bias problem.

A number of authors have considered modified training procedures for greedy models, or for locally normalized models. Daumé III et al. (2009) introduce Searn, an algorithm that allows a classifier making greedy decisions to become more robust to errors made in previous decisions. Goldberg and Nivre (2013) describe improvements to a greedy parsing approach that makes use of methods from imitation learning (Ross et al., 2011) to augment the training set. Note that these methods are focused on greedy models: they are unlikely to solve the label bias problem when used in conjunction with beam search, given that the problem is one of expressivity of the underlying model. More recent work (Yazdani and Henderson, 2015; Vaswani and Sagae, 2016) has augmented locally normalized models with *correctness probabilities* or *error states*, effectively adding a step after every decision where the probability of correctness of the resulting structure is evaluated. This gives considerable gains over a locally normalized model, although performance is lower than our full globally normalized approach.

### 6 Conclusions

We presented a simple and yet powerful model architecture that produces state-of-the-art results for POS tagging, dependency parsing and sentence compression. Our model combines the flexibility of transition-based algorithms and the modeling power of neural networks. Our results demon-

---

[3]This setting may be important in some applications, where for example parse structures for sentence prefixes are required, or where the input is received one word at a time and online processing is beneficial.

strate that feed-forward network without recurrence can outperform recurrent models such as LSTMs when they are trained with global normalization. We further support our empirical findings with a proof showing that global normalization helps the model overcome the label bias problem from which locally normalized models suffer.

## Acknowledgements

## References

[Alberti et al.2015] Chris Alberti, David Weiss, Greg Coppola, and Slav Petrov. 2015. Improved transition-based parsing and tagging with neural networks. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1354–1359.

[Ballesteros et al.2015] Miguel Ballesteros, Chris Dyer, and Noah A. Smith. 2015. Improved transition-based parsing by modeling characters instead of words with LSTMs. pages 349–359.

[Bohnet and Nivre2012] Bernd Bohnet and Joakim Nivre. 2012. A transition-based system for joint part-of-speech tagging and labeled non-projective dependency parsing. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 1455–1465.

[Bottou and LeCun2005] Léon Bottou and Yann Le-Cun. 2005. Graph transformer networks for image recognition. *Bulletin of the International Statistical Institute (ISI)*.

[Bottou et al.1997] Léon Bottou, Yann Le Cun, and Yoshua Bengio. 1997. Global training of document processing systems using graph transformer networks. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 489–493.

[Bottou1991] Léon Bottou. 1991. *Une approche théorique de lapprentissage connexionniste: Applications à la reconnaissance de la parole*. Ph.D. thesis, Doctoral dissertation, Universite de Paris XI.

[Chen and Manning2014] Danqi Chen and Christopher D. Manning. 2014. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 740–750.

[Collins and Roark2004] Michael Collins and Brian Roark. 2004. Incremental parsing with the perceptron algorithm. In *Proceedings of the 42nd Meeting of the Association for Computational Linguistics (ACL'04)*, pages 111–118.

[Daumé III et al.2009] Hal Daumé III, John Langford, and Daniel Marcu. 2009. Search-based structured prediction. *Machine Learning Journal (MLJ)*.

[De Marneffe et al.2006] Marie-Catherine De Marneffe, Bill MacCartney, and Christopher D. Manning. 2006. Generating typed dependency parses from phrase structure parses. In *Proceedings of Fifth International Conference on Language Resources and Evaluation*, pages 449–454.

[Do and Artires2010] Trinh Minh Tri Do and Thierry Artires. 2010. Neural conditional random fields. In *International Conference on Artificial Intelligence and Statistics*, volume 9, pages 177–184.

[Durrett and Klein2015] Greg Durrett and Dan Klein. 2015. Neural crf parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, pages 302–312.

[Dyer et al.2015] Chris Dyer, Miguel Ballesteros, Wang Ling, Austin Matthews, and Noah A. Smith. 2015. Transition-based dependency parsing with stack long short-term memory. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*, pages 334–343.

[Filippova et al.2015] Katja Filippova, Enrique Alfonseca, Carlos A. Colmenares, Łukasz Kaiser, and Oriol Vinyals. 2015. Sentence compression by deletion with lstms. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 360–368.

[Goldberg and Nivre2013] Yoav Goldberg and Joakim Nivre. 2013. Training deterministic parsers with non-deterministic oracles. *Transactions of the Association for Computational Linguistics*, 1:403–414.

[Hajič et al.2009] Jan Hajič, Massimiliano Ciaramita, Richard Johansson, Daisuke Kawahara, Maria Antònia Martí, Lluís Màrquez, Adam Meyers, Joakim Nivre, Sebastian Padó, Jan Štěpánek, Pavel Straňák, Mihai Surdeanu, Nianwen Xue, and Yi Zhang. 2009. The conll-2009 shared task: Syntactic and semantic dependencies in multiple languages. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning: Shared Task*, pages 1–18.

[Henderson2003] James Henderson. 2003. Inducing history representations for broad coverage statistical parsing. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 24–31.

[Henderson2004] James Henderson. 2004. Discriminative training of a neural network statistical parser. In *Proceedings of the 42nd Meeting of the Association for Computational Linguistics (ACL'04)*, pages 95–102.

[Hochreiter and Schmidhuber1997] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

[Hovy et al.2006] Eduard Hovy, Mitchell Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2006. Ontonotes: The 90% solution. In *Proceedings of the Human Language Technology Conference of the NAACL, Short Papers*, pages 57–60.

[Huang et al.2015] Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional LSTM-CRF models for sequence tagging. *arXiv preprint arXiv:1508.01991*.

[Judge et al.2006] John Judge, Aoife Cahill, and Josef van Genabith. 2006. Questionbank: Creating a corpus of parse-annotated questions. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 497–504.

[Lafferty et al.2001] John Lafferty, Andrew McCallum, and Fernando Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 282–289.

[Le Cun et al.1998] Yann Le Cun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient based learning applied to document recognition. *Proceedings of IEEE*, 86(11):2278–2324.

[Lei et al.2014] Tao Lei, Yu Xin, Yuan Zhang, Regina Barzilay, and Tommi Jaakkola. 2014. Low-rank tensors for scoring dependency structures. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pages 1381–1391.

[Liang et al.2008] Percy Liang, Hal Daumé, III, and Dan Klein. 2008. Structure compilation: Trading structure for features. In *Proceedings of the 25th International Conference on Machine Learning*, pages 592–599.

[Ling et al.2015] Wang Ling, Chris Dyer, Alan W Black, Isabel Trancoso, Ramon Fermandez, Silvio Amir, Luis Marujo, and Tiago Luis. 2015. Finding function in form: Compositional character models for open vocabulary word representation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1520–1530.

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

[Martins et al.2013] Andre Martins, Miguel Almeida, and Noah A. Smith. 2013. Turning on the turbo: Fast third-order non-projective turbo parsers. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, pages 617–622.

[Nivre2006] Joakim Nivre. 2006. *Inductive Dependency Parsing*. Springer-Verlag New York, Inc.

[Nivre2009] Joakim Nivre. 2009. Non-projective dependency parsing in expected linear time. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 351–359.

[Peng et al.2009] Jian Peng, Liefeng Bo, and Jinbo Xu. 2009. Conditional neural fields. In *Advances in Neural Information Processing Systems 22*, pages 1419–1427.

[Petrov and McDonald2012] Slav Petrov and Ryan McDonald. 2012. Overview of the 2012 shared task on parsing the web. Notes of the First Workshop on Syntactic Analysis of Non-Canonical Language (SANCL).

[Ross et al.2011] Stéphane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. 2011. No-regret reductions for imitation learning and structured prediction. *AISTATS*.

[Vaswani and Sagae2016] Ashish Vaswani and Kenji Sagae. 2016. Efficient structured inference for transition-based parsing with neural networks and error states. *Transactions of the Association for Computational Linguistics*, to appear.

[Vinyals et al.2015] Oriol Vinyals, Łukasz Kaiser, Terry Koo, Slav Petrov, Ilya Sutskever, and Geoffrey Hinton. 2015. Grammar as a foreign language. In *Advances in Neural Information Processing Systems 28*, pages 2755–2763.

[Watanabe and Sumita2015] Taro Watanabe and Eiichiro Sumita. 2015. Transition-based neural constituent parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, pages 1169–1179.

[Weiss et al.2015] David Weiss, Chris Alberti, Michael Collins, and Slav Petrov. 2015. Structured training for neural network transition-based parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*, pages 323–333.

[Yao et al.2014] Kaisheng Yao, Baolin Peng, Geoffrey Zweig, Dong Yu, Xiaolong Li, and Feng Gao. 2014. Recurrent conditional random field for language understanding. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '14)*.

[Yazdani and Henderson2015] Majid Yazdani and James Henderson. 2015. Incremental recurrent neural network dependency parser with search-based discriminative training. In *Proceedings of the Nineteenth Conference on Computational Natural Language Learning*, pages 142–152.

[Zhang and McDonald2014] Hao Zhang and Ryan McDonald. 2014. Enforcing structural diversity in cube-pruned dependency parsing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pages 656–661.

[Zheng et al.2015] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip H. S. Torr. 2015. Conditional random fields as recurrent neural networks. In *The IEEE International Conference on Computer Vision (ICCV)*.

[Zhou and Xu2015] Jie Zhou and Wei Xu. 2015. End-to-end learning of semantic role labeling using recurrent neural networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, pages 1127–1137.

[Zhou et al.2015] Hao Zhou, Yue Zhang, and Jiajun Chen. 2015. A neural probabilistic structured-prediction model for transition-based dependency parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*, pages 1213–1222.