

# Encoding Sentences with Graph Convolutional Networks for Semantic Role Labeling

Diego Marcheggiani<sup>1</sup>

Ivan Titov<sup>1,2</sup>

<sup>1</sup>ILLC, University of Amsterdam

<sup>2</sup>ILCC, School of Informatics, University of Edinburgh

marcheggiani@uva.nl

ititov@inf.ed.ac.uk

## Abstract

Semantic role labeling (SRL) is the task of identifying the predicate-argument structure of a sentence. It is typically regarded as an important step in the standard NLP pipeline. As the semantic representations are closely related to syntactic ones, we exploit syntactic information in our model. We propose a version of graph convolutional networks (GCNs), a recent class of neural networks operating on graphs, suited to model syntactic dependency graphs. GCNs over syntactic dependency trees are used as sentence encoders, producing latent feature representations of words in a sentence. We observe that GCN layers are complementary to LSTM ones: when we stack both GCN and LSTM layers, we obtain a substantial improvement over an already state-of-the-art LSTM SRL model, resulting in the best reported scores on the standard benchmark (CoNLL-2009) both for Chinese and English.

## 1 Introduction

Semantic role labeling (SRL) (Gildea and Jurafsky, 2002) can be informally described as the task of discovering *who* did *what* to *whom*. For example, consider an SRL dependency graph shown above the sentence in Figure 1. Formally, the task includes (1) detection of predicates (e.g., *makes*); (2) labeling the predicates with a sense from a sense inventory (e.g., *make.01*); (3) identifying and assigning arguments to *semantic roles* (e.g., *Sequa* is A0, i.e., an agent / ‘doer’ for the corresponding predicate, and *engines* is A1, i.e., a patient / ‘an affected entity’). SRL is often regarded

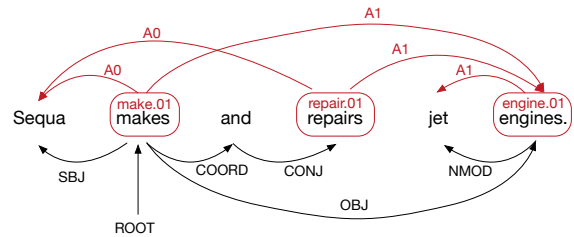


Figure 1: An example sentence annotated with semantic (top) and syntactic dependencies (bottom).

as an important step in the standard NLP pipeline, providing information to downstream tasks such as information extraction and question answering.

The semantic representations are closely related to syntactic ones, even though the syntax-semantics interface is far from trivial (Levin, 1993). For example, one can observe that many arcs in the syntactic dependency graph (shown in black below the sentence in Figure 1) are mirrored in the semantic dependency graph. **Given these similarities and also because of availability of accurate syntactic parsers for many languages, it seems natural to exploit syntactic information when predicting semantics.** Though historically most SRL approaches did rely on syntax (Thompson et al., 2003; Pradhan et al., 2005; Punyakanok et al., 2008; Johansson and Nugues, 2008), the last generation of SRL models put syntax aside in favor of neural sequence models, namely LSTMs (Zhou and Xu, 2015; Marcheggiani et al., 2017), and outperformed syntactically-driven methods on standard benchmarks. **We believe that one of the reasons for this radical choice is the lack of simple and effective methods for incorporating syntactic information into sequential neural networks (namely, at the level of words).** In this paper we

propose one way how to address this limitation.

Specifically, we rely on graph convolutional networks (GCNs) (Duvenaud et al., 2015; Kipf and Welling, 2017; Kearnes et al., 2016), a recent class of multilayer neural networks operating on graphs. For every node in the graph (in our case a word in a sentence), GCN encodes relevant information about its neighborhood as a real-valued feature vector. GCNs have been studied largely in the context of undirected unlabeled graphs. We introduce a version of **GCNs for modeling syntactic dependency structures and generally applicable to labeled directed graphs**.

One layer GCN encodes only information about immediate neighbors and  $K$  layers are needed to encode  $K$ -order neighborhoods (i.e., information about nodes at most  $K$  hops away). This contrasts with recurrent and recursive neural networks (Elman, 1990; Socher et al., 2013) which, at least in theory, can capture statistical dependencies across unbounded paths in a trees or in a sequence. However, as we will further discuss in Section 3.3, this is not a serious limitation when GCNs are used in combination with encoders based on recurrent networks (LSTMs). When we stack GCNs on top of LSTM layers, we obtain a substantial improvement over an already state-of-the-art LSTM SRL model, resulting in the best reported scores on the standard benchmark (CoNLL-2009), both for English and Chinese.<sup>1</sup>

Interestingly, again unlike recursive neural networks, GCNs do not constrain the graph to be a tree. We believe that there are many applications in NLP, where GCN-based encoders of sentences or even documents can be used to incorporate knowledge about linguistic structures (e.g., representations of syntax, semantics or discourse). For example, GCNs can take as input combined syntactic-semantic graphs (e.g., the entire graph from Figure 1) and be used within downstream tasks such as machine translation or question answering. However, we leave this for future work and here solely focus on SRL.

The contributions of this paper can be summarized as follows:

- we are the first to show that GCNs are effective for NLP;
- we propose a generalization of GCNs suited

to encoding syntactic information at word level;

- we propose a GCN-based SRL model and obtain state-of-the-art results on English and Chinese portions of the CoNLL-2009 dataset;
- we show that bidirectional LSTMs and syntax-based GCNs have complementary modeling power.

## 2 Graph Convolutional Networks

In this section we describe GCNs of Kipf and Welling (2017). Please refer to Gilmer et al. (2017) for a comprehensive overview of GCN versions.

GCNs are neural networks operating on graphs and inducing features of nodes (i.e., real-valued vectors / embeddings) based on properties of their neighborhoods. In Kipf and Welling (2017), they were shown to be very effective for the node classification task: the classifier was estimated jointly with a GCN, so that the induced node features were informative for the node classification problem. Depending on how many layers of convolution are used, GCNs can capture information only about immediate neighbors (with one layer of convolution) or any nodes at most  $K$  hops away (if  $K$  layers are stacked on top of each other).

More formally, consider an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  ( $|\mathcal{V}| = n$ ) and  $\mathcal{E}$  are sets of nodes and edges, respectively. Kipf and Welling (2017) assume that edges contain all the self-loops, i.e.,  $(v, v) \in \mathcal{E}$  for any  $v$ . We can define a matrix  $X \in \mathbb{R}^{m \times n}$  with each its column  $x_v \in \mathbb{R}^m$  ( $v \in \mathcal{V}$ ) encoding node features. The vectors can either encode genuine features (e.g., this vector can encode the title of a paper if citation graphs are considered) or be a one-hot vector. The node representation, encoding information about its immediate neighbors, is computed as

$$h_v = \text{ReLU} \left( \sum_{u \in \mathcal{N}(v)} (W x_u + b) \right), \quad (1)$$

where  $W \in \mathbb{R}^{m \times m}$  and  $b \in \mathbb{R}^m$  are a weight matrix and a bias, respectively;  $\mathcal{N}(v)$  are neighbors of  $v$ ;  $\text{ReLU}$  is the rectifier linear unit activation function.<sup>2</sup> Note that  $v \in \mathcal{N}(v)$  (because of self-loops), so the input feature representation of  $v$  (i.e.  $x_v$ ) affects its induced representation  $h_v$ .

<sup>1</sup>The code is available at <https://github.com/diegma/neural-dep-srl>.

<sup>2</sup>We dropped normalization factors used in Kipf and Welling (2017), as they are not used in our syntactic GCNs.

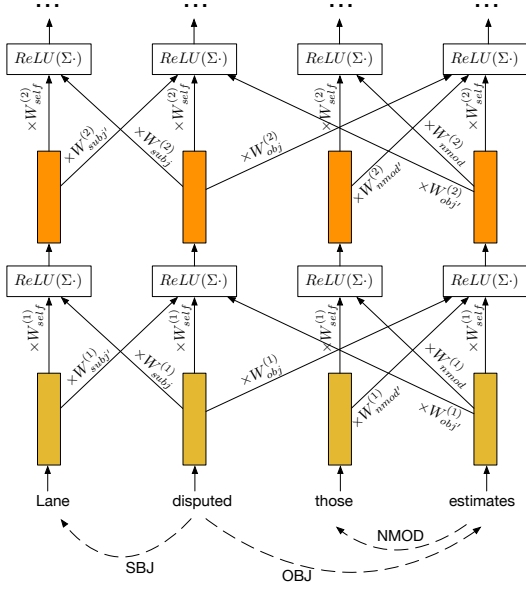


Figure 2: A simplified syntactic GCN (bias terms and gates are omitted); the syntactic graph of the sentence is shown with dashed lines at the bottom. Parameter matrices are sub-indexed with syntactic functions, and apostrophes (e.g., *subj'*) signify that information flows in the direction opposite of the dependency arcs (i.e., from dependents to heads).

As in standard convolutional networks (LeCun et al., 2001), by stacking GCN layers one can incorporate higher degree neighborhoods:

$$h_v^{(k+1)} = \text{ReLU} \left( \sum_{u \in \mathcal{N}(v)} W^{(k)} h_u^{(k)} + b^{(k)} \right)$$

where  $k$  denotes the layer number and  $h_v^{(1)} = x_v$ .

### 3 Syntactic GCNs

As syntactic dependency trees are directed and labeled (we refer to the dependency labels as *syntactic functions*), we first need to modify the computation in order to incorporate label information (Section 3.1). In the subsequent section, we incorporate gates in GCNs, so that the model can decide which edges are more relevant to the task in question. Having gates is also important as we rely on automatically predicted syntactic representations, and the gates can detect and downweight potentially erroneous edges.

#### 3.1 Incorporating directions and labels

Now, we introduce a generalization of GCNs appropriate for syntactic dependency trees, and in

general, for directed labeled graphs. First note that there is no reason to assume that information flows only along the syntactic dependency arcs (e.g., from *makes* to *Sequa*), so we allow it to flow in the opposite direction as well (i.e., from dependents to heads). We use a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where the edge set contains all pairs of nodes (i.e., words) adjacent in the dependency tree. In our example, both  $(\text{Sequa}, \text{makes})$  and  $(\text{makes}, \text{Sequa})$  belong to the edge set. The graph is labeled, and the label  $L(u, v)$  for  $(u, v) \in \mathcal{E}$  contains both information about the syntactic function and indicates whether the edge is in the same or opposite direction as the syntactic dependency arc. For example, the label for  $(\text{makes}, \text{Sequa})$  is *subj*, whereas the label for  $(\text{Sequa}, \text{makes})$  is *subj'*, with the apostrophe indicating that the edge is in the direction opposite to the corresponding syntactic arc. Similarly, self-loops will have label *self*. Consequently, we can simply assume that the GCN parameters are label-specific, resulting in the following computation, also illustrated in Figure 2:

$$h_v^{(k+1)} = \text{ReLU} \left( \sum_{u \in \mathcal{N}(v)} W_{L(u,v)}^{(k)} h_u^{(k)} + b_{L(u,v)}^{(k)} \right).$$

This model is over-parameterized,<sup>3</sup> especially given that SRL datasets are moderately sized, by deep learning standards. So instead of learning the GCN parameters directly, we define them as

$$W_{L(u,v)}^{(k)} = V_{\text{dir}(u,v)}^{(k)}, \quad (2)$$

where  $\text{dir}(u, v)$  indicates whether the edge  $(u, v)$  is directed (1) along, (2) in the opposite direction to the syntactic dependency arc, or (3) is a self-loop;  $V_{\text{dir}(u,v)}^{(k)} \in \mathbb{R}^{m \times m}$ . **Our simplification captures the intuition that information should be propagated differently along edges depending whether this is a head-to-dependent or dependent-to-head edge (i.e., along or opposite the corresponding syntactic arc) and whether it is a self-loop.** So we do not share any parameters between these three very different edge types. Syntactic functions are important, but perhaps less crucial, so they are encoded only in the feature vectors  $b_{L(u,v)}$ .

#### 3.2 Edge-wise gating

Uniformly accepting information from all neighboring nodes may not be appropriate for the SRL

<sup>3</sup>Chinese and English CoNLL-2009 datasets used 41 and 48 different syntactic functions, which would result in having 83 and 97 different matrices in every layer, respectively.

setting. For example, we see in Figure 1 that many semantic arcs just mirror their syntactic counterparts, so they may need to be up-weighted. Moreover, we rely on automatically predicted syntactic structures, and, even for English, syntactic parsers are far from being perfect, especially when used out-of-domain. **It is risky for a downstream application to rely on a potentially wrong syntactic edge, so the corresponding message in the neural network may need to be down-weighted.**

In order to address the above issues, inspired by recent literature (van den Oord et al., 2016; Dauphin et al., 2016), we calculate for each edge node pair a scalar gate of the form

$$g_{u,v}^{(k)} = \sigma \left( h_u^{(k)} \cdot \hat{v}_{dir(u,v)}^{(k)} + \hat{b}_{L(u,v)}^{(k)} \right), \quad (3)$$

where  $\sigma$  is the logistic sigmoid function,  $\hat{v}_{dir(u,v)}^{(k)} \in \mathbb{R}^m$  and  $\hat{b}_{L(u,v)}^{(k)} \in \mathbb{R}$  are weights and a bias for the gate. With this additional gating mechanism, the final syntactic GCN computation is formulated as

$$h_v^{(k+1)} = ReLU \left( \sum_{u \in \mathcal{N}(v)} g_{v,u}^{(k)} (V_{dir(u,v)}^{(k)} h_u^{(k)} + b_{L(u,v)}^{(k)}) \right). \quad (4)$$

### 3.3 Complementarity of GCNs and LSTMs

The inability of GCNs to capture dependencies between nodes far away from each other in the graph may seem like a serious problem, especially in the context of SRL: paths between predicates and arguments often include many dependency arcs (Roth and Lapata, 2016). However, when graph convolution is performed on top of LSTM states (i.e., LSTM states serve as input  $x_v = h_v^{(1)}$  to GCN) rather than static word embeddings, GCN may not need to capture more than a couple of hops.

To elaborate on this, let us speculate what role GCNs would play when used in combinations with LSTMs, given that LSTMs have already been shown very effective for SRL (Zhou and Xu, 2015; Marcheggiani et al., 2017). Though LSTMs are capable of capturing at least some degree of syntax (Linzen et al., 2016) without explicit syntactic supervision, SRL datasets are moderately sized, so LSTM models may still struggle with harder cases. Typically, harder cases for SRL involve arguments far away from their predicates. In fact, 20% and 30% of arguments are more than 5 tokens away from their predicate, in our English and

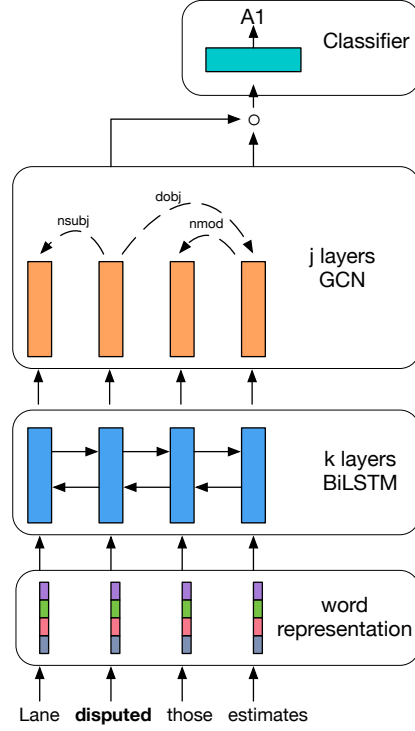


Figure 3: Predicting an argument and its label with an LSTM + GCN encoder.

Chinese collections, respectively. However, if we imagine that we can ‘teleport’ even over a single (longest) syntactic dependency edge, the ‘distance’ would shrink: only 9% and 13% arguments will now be more than 5 LSTM steps away (again for English and Chinese, respectively). GCNs provide this ‘teleportation’ capability. These observations suggest that LSTMs and GCNs may be complementary, and we will see that empirical results support this intuition.

## 4 Syntax-Aware Neural SRL Encoder

In this work, we build our semantic role labeler on top of the syntax-agnostic LSTM-based SRL model of Marcheggiani et al. (2017), which already achieves state-of-the-art results on the CoNLL-2009 English dataset. Following their approach we employ the same bidirectional (BiLSTM) encoder and enrich it with a syntactic GCN.

The CoNLL-2009 benchmark assumes that predicate positions are already marked in the test set (e.g., we would know that *makes*, *repairs* and *engines* in Figure 1 are predicates), so no predicate identification is needed. Also, as we focus here solely on identifying arguments and labeling them with semantic roles, for predicate disambiguation



(i.e., marking *makes* as *make.01*) we use of an off-the-shelf disambiguation model (Roth and Lapata, 2016; Björkelund et al., 2009). As in Marcheggiani et al. (2017) and in most previous work, we process individual predicates in isolation, so for each predicate, our task reduces to a sequence labeling problem. That is, given a predicate (e.g., *disputed* in Figure 3) one needs to identify and label all its arguments (e.g., label *estimates* as A1 and label *those* as ‘NULL’, indicating that *those* is not an argument of *disputed*).

The semantic role labeler we propose is composed of four components (see Figure 3):

- look-ups of word embeddings;
- a BiLSTM encoder that takes as input the word representation of each word in a sentence;
- a syntax-based GCN encoder that re-encodes the BiLSTM representation based on the automatically predicted syntactic structure of the sentence;
- a role classifier that takes as input the GCN representation of the candidate argument and the representation of the predicate to predict the role associated with the candidate word.

#### 4.1 Word representations

For each word  $w_i$  in the considered sentence, we create a sentence-specific word representation  $x_i$ . We represent each word  $w$  as the concatenation of four vectors:<sup>4</sup> a randomly initialized word embedding  $x^{re} \in \mathbb{R}^{d_w}$ , a pre-trained word embedding  $x^{pe} \in \mathbb{R}^{d_w}$  estimated on an external text collection, a randomly initialized part-of-speech tag embedding  $x^{pos} \in \mathbb{R}^{d_p}$  and a randomly initialized lemma embedding  $x^{le} \in \mathbb{R}^{d_l}$  (active only if the word is a predicate). The randomly initialized embeddings  $x^{re}$ ,  $x^{pos}$ , and  $x^{le}$  are fine-tuned during training, while the pre-trained ones are kept fixed. The final word representation is given by  $x = x^{re} \circ x^{pe} \circ x^{pos} \circ x^{le}$ , where  $\circ$  represents the concatenation operator.

#### 4.2 Bidirectional LSTM layer

One of the most popular and effective ways to represent sequences, such as sentences (Mikolov et al., 2010), is to use recurrent neural networks

(RNN) (Elman, 1990). In particular their gated versions, Long Short-Term Memory (LSTM) networks (Hochreiter and Schmidhuber, 1997) and Gated Recurrent Units (GRU) (Cho et al., 2014), have proven effective in modeling long sequences (Chiu and Nichols, 2016; Sutskever et al., 2014).

Formally, an LSTM can be defined as a function  $LSTM_\theta(x_{1:i})$  that takes as input the sequence  $x_{1:i}$  and returns a hidden state  $h_i \in \mathbb{R}^{d_h}$ . This state can be regarded as a representation of the sentence from the start to the position  $i$ , or, in other words, it encodes the word at position  $i$  along with its left context. However, the right context is also important, so Bidirectional LSTMs (Graves, 2008) use two LSTMs: one for the forward pass, and another for the backward pass,  $LSTM_F$  and  $LSTM_B$ , respectively. By concatenating the states of both LSTMs, we create a complete context-aware representation of a word  $BiLSTM(x_{1:n}, i) = LSTM_F(x_{1:i}) \circ LSTM_B(x_{n:i})$ . We follow Marcheggiani et al. (2017) and stack  $J$  layers of bidirectional LSTMs, where each layer takes the lower layer as its input.

#### 4.3 Graph convolutional layer

The representation calculated with the BiLSTM encoder is fed as input to a GCN of the form defined in Equation (4). The neighboring nodes of a node  $v$ , namely  $\mathcal{N}(v)$ , and their relations to  $v$  are predicted by an external syntactic parser.

#### 4.4 Semantic role classifier

The classifier predicts semantic roles of words given the predicate while relying on word representations provided by GCN; we concatenate hidden states of the candidate argument word and the predicate word and use them as input to a classifier (Figure 3, top). The softmax classifier computes the probability of the role (including special ‘NULL’ role):

$$p(r|t_i, t_p, l) \propto \exp(W_{l,r}(t_i \circ t_p)), \quad (5)$$

where  $t_i$  and  $t_p$  are representations produced by the graph convolutional encoder,  $l$  is the lemma of predicate  $p$ , and the symbol  $\propto$  signifies proportionality.<sup>5</sup> As FitzGerald et al. (2015) and Marcheggiani et al. (2017), instead of using a fixed matrix  $W_{l,r}$  or simply assuming that  $W_{l,r} = W_r$ ,

<sup>4</sup>We drop the index  $i$  from the notation for the sake of brevity.

<sup>5</sup>We abuse the notation and refer as  $p$  both to the predicate word and to its position in the sentence.

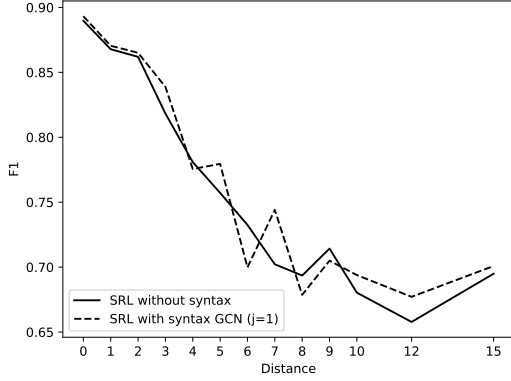


Figure 4:  $F_1$  as function of word distance. The distance starts from zero, since nominal predicates can be arguments of themselves.

we jointly embed the role  $r$  and predicate lemma  $l$  using a non-linear transformation:

$$W_{l,r} = \text{ReLU}(U(q_l \circ q_r)), \quad (6)$$

where  $U$  is a parameter matrix, whereas  $q_l \in \mathbb{R}^{d_l}$  and  $q_r \in \mathbb{R}^{d_r}$  are randomly initialized embeddings of predicate lemmas and roles. **In this way each role prediction is predicate-specific**, and, at the same time, we expect to learn a good representation for roles associated with infrequent predicates. As our training objective we use the categorical cross-entropy.

## 5 Experiments

### 5.1 Datasets and parameters

We tested the proposed SRL model on the English and Chinese CoNLL-2009 dataset with standard splits into training, test and development sets. The predicted POS tags for both languages were provided by the CoNLL-2009 shared-task organizers. For the predicate disambiguator we used the ones from [Roth and Lapata \(2016\)](#) for English and from [Björkelund et al. \(2009\)](#) for Chinese. We parsed English sentences with the BIST Parser ([Kiperwasser and Goldberg, 2016](#)), whereas for Chinese we used automatically predicted parses provided by the CoNLL-2009 shared-task organizers.

For English, we used external embeddings of [Dyer et al. \(2015\)](#), learned using the structured skip n-gram approach of [Ling et al. \(2015\)](#). For Chinese we used external embeddings produced with the neural language model of [Bengio et al. \(2003\)](#). We used *edge dropout* in GCN: when

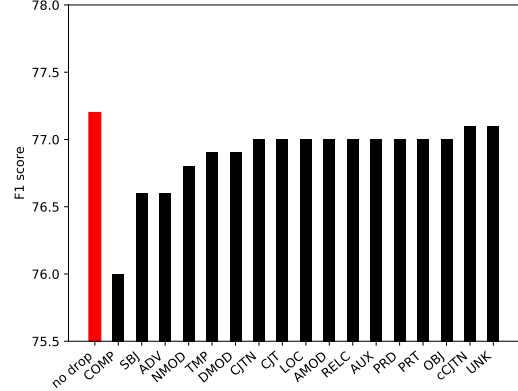


Figure 5: Performance with dependency arcs of given type dropped, on Chinese development set.

| System (English)                 | P    | R    | $F_1$ |
|----------------------------------|------|------|-------|
| LSTMs                            | 84.3 | 81.1 | 82.7  |
| LSTMs + GCNs ( $K=1$ )           | 85.2 | 81.6 | 83.3  |
| LSTMs + GCNs ( $K=2$ )           | 84.1 | 81.4 | 82.7  |
| LSTMs + GCNs ( $K=1$ ), no gates | 84.7 | 81.4 | 83.0  |
| GCNs (no LSTMs), $K=1$           | 79.9 | 70.4 | 74.9  |
| GCNs (no LSTMs), $K=2$           | 83.4 | 74.6 | 78.7  |
| GCNs (no LSTMs), $K=3$           | 83.6 | 75.8 | 79.5  |
| GCNs (no LSTMs), $K=4$           | 82.7 | 76.0 | 79.2  |

Table 1: SRL results without predicate disambiguation on the English development set.

computing  $h_v^{(k)}$ , we ignore each node  $v \in \mathcal{N}(v)$  with probability  $\beta$ . Adam ([Kingma and Ba, 2015](#)) was used as an optimizer. The hyperparameter tuning and all model selection were performed on the English development set; the chosen values are shown in Appendix.

### 5.2 Results and discussion

In order to show that GCN layers are effective, we first compare our model against its version which lacks GCN layers (i.e. essentially the model of [Marcheggiani et al. \(2017\)](#)). Importantly, to measure the genuine contribution of GCNs, we first tuned this syntax-agnostic model (e.g., the number of LSTM layers) to get best possible performance on the development set.<sup>6</sup>

We compare the syntax-agnostic model with 3 syntax-aware versions: one GCN layer over syntax ( $K = 1$ ), one layer GCN without gates and two GCN layers ( $K = 2$ ). As we rely on the same

<sup>6</sup>For example, if we would have used only one layer of LSTMs, gains from using GCNs would be even larger.

| System (Chinese)             | P    | R    | F <sub>1</sub> |
|------------------------------|------|------|----------------|
| LSTMs                        | 78.3 | 72.3 | 75.2           |
| LSTMs + GCNs (K=1)           | 79.9 | 74.4 | 77.1           |
| LSTMs + GCNs (K=2)           | 78.7 | 74.0 | 76.2           |
| LSTMs + GCNs (K=1), no gates | 78.2 | 74.8 | 76.5           |
| GCNs (no LSTMs), K=1         | 78.7 | 58.5 | 67.1           |
| GCNs (no LSTMs), K=2         | 79.7 | 62.7 | 70.1           |
| GCNs (no LSTMs), K=3         | 76.8 | 66.8 | 71.4           |
| GCNs (no LSTMs), K=4         | 79.1 | 63.5 | 70.4           |

Table 2: SRL results without predicate disambiguation on the Chinese development set.

| System                              | P           | R           | F <sub>1</sub> |
|-------------------------------------|-------------|-------------|----------------|
| Lei et al. (2015) (local)           | -           | -           | 86.6           |
| FitzGerald et al. (2015) (local)    | -           | -           | 86.7           |
| Roth and Lapata (2016) (local)      | 88.1        | 85.3        | 86.7           |
| Marcheggiani et al. (2017) (local)  | 88.6        | 86.7        | 87.6           |
| <b>Ours (local)</b>                 | <b>89.1</b> | <b>86.8</b> | <b>88.0</b>    |
| Björkelund et al. (2010) (global)   | 88.6        | 85.2        | 86.9           |
| FitzGerald et al. (2015) (global)   | -           | -           | 87.3           |
| Foland and Martin (2015) (global)   | -           | -           | 86.0           |
| Swayamdipta et al. (2016) (global)  | -           | -           | 85.0           |
| Roth and Lapata (2016) (global)     | 90.0        | 85.5        | 87.7           |
| FitzGerald et al. (2015) (ensemble) | -           | -           | 87.7           |
| Roth and Lapata (2016) (ensemble)   | 90.3        | 85.7        | 87.9           |
| <b>Ours (ensemble 3x)</b>           | <b>90.5</b> | <b>87.7</b> | <b>89.1</b>    |

Table 3: Results on the test set for English.

off-the-shelf disambiguator for all versions of the model, in Table 1 and 2 we report SRL-only scores (i.e., predicate disambiguation is not evaluated) on the English and Chinese development sets. For both datasets, the syntax-aware model with one GCN layers ( $K = 1$ ) performs the best, outperforming the LSTM version by 1.9% and 0.6% for Chinese and English, respectively. The reasons why the improvements on Chinese are much larger are not entirely clear (e.g., both languages are relative fixed word order ones, and the syntactic parses for Chinese are considerably less accurate), this may be attributed to a higher proportion of long-distance dependencies between predicates and arguments in Chinese (see Section 3.3). Edge-wise gating (Section 3.2) also appears important: removing gates leads to a drop of 0.3% F<sub>1</sub> for English and 0.6% F<sub>1</sub> for Chinese.

Stacking two GCN layers does not give any benefit. When BiLSTM layers are dropped altogether, stacking two layers ( $K = 2$ ) of GCNs greatly improves the performance, resulting in a 3.8% jump in F<sub>1</sub> for English and a 3.0% jump in F<sub>1</sub> for Chi-

| System                            | P           | R           | F <sub>1</sub> |
|-----------------------------------|-------------|-------------|----------------|
| Zhao et al. (2009) (global)       | 80.4        | 75.2        | 77.7           |
| Björkelund et al. (2009) (global) | 82.4        | 75.1        | 78.6           |
| Roth and Lapata (2016) (global)   | 83.2        | 75.9        | 79.4           |
| <b>Ours (local)</b>               | <b>84.6</b> | <b>80.4</b> | <b>82.5</b>    |

Table 4: Results on the Chinese test set.

nese. Adding a 3rd layer of GCN ( $K = 3$ ) further improves the performance.<sup>7</sup> This suggests that extra GCN layers are effective but largely redundant with respect to what LSTMs already capture.

In Figure 4, we show the  $F_1$  scores results on the English development set as a function of the distance, in terms of tokens, between a candidate argument and its predicate. As expected, GCNs appear to be more beneficial for long distance dependencies, as shorter ones are already accurately captured by the LSTM encoder.

We looked closer in contribution of specific dependency relations for Chinese. In order to assess this without retraining the model multiple times, we drop all dependencies of a given type at test time (one type at a time, only for types appearing over 300 times in the development set) and observe changes in performance. In Figure 5, we see that the most informative dependency is COMP (complement). Relative clauses in Chinese are very frequent and typically marked with particle 的 (de). The relative clause will syntactically depend on 的 as COMP, so COMP encodes important information about predicate-argument structure. These are often long-distance dependencies and may not be accurately captured by LSTMs. Although TMP (temporal) dependencies are not as frequent ( $\sim 2\%$  of all dependencies), they are also important: temporal information is mirrored in semantic roles.

In order to compare to previous work, in Table 3 we report test results on the English in-domain (WSJ) evaluation data. Our model is *local*, as all the argument detection and labeling decisions are conditionally independent: their interaction is captured solely by the LSTM+GCN encoder. This makes our model fast and simple, though, as shown in previous work, *global* modeling of the structured output is beneficial.<sup>8</sup> We leave this extension for future work. Interestingly,

<sup>7</sup>Note that GCN layers are computationally cheaper than LSTM ones, even in our non-optimized implementation.

<sup>8</sup>As seen in Table 3, labelers of FitzGerald et al. (2015) and Roth and Lapata (2016) gained 0.6-1.0%.

| System                              | P           | R           | F <sub>1</sub> |
|-------------------------------------|-------------|-------------|----------------|
| Lei et al. (2015) (local)           | -           | -           | 75.6           |
| FitzGerald et al. (2015) (local)    | -           | -           | 75.2           |
| Roth and Lapata (2016) (local)      | 76.9        | 73.8        | 75.3           |
| Marcheggiani et al. (2017) (local)  | 78.9        | 75.7        | 77.3           |
| <b>Ours (local)</b>                 | <b>78.5</b> | <b>75.9</b> | <b>77.2</b>    |
| Björkelund et al. (2010) (global)   | 77.9        | 73.6        | 75.7           |
| FitzGerald et al. (2015) (global)   | -           | -           | 75.2           |
| Foland and Martin (2015) (global)   | -           | -           | 75.9           |
| Roth and Lapata (2016) (global)     | 78.6        | 73.8        | 76.1           |
| FitzGerald et al. (2015) (ensemble) | -           | -           | 75.5           |
| Roth and Lapata (2016) (ensemble)   | 79.7        | 73.6        | 76.5           |
| <b>Ours (ensemble 3x)</b>           | <b>80.8</b> | <b>77.1</b> | <b>78.9</b>    |

Table 5: Results on the out-of-domain test set.

we outperform even the best global model and the best ensemble of global models, without using global modeling or ensembles. When we create an ensemble of 3 models with the product-of-expert combination rule, we improve by 1.2% over the best previous result, achieving 89.1% F<sub>1</sub>.<sup>9</sup>

For Chinese (Table 4), our best model outperforms the state-of-the-art model of Roth and Lapata (2016) by even larger margin of 3.1%.

For English, in the CoNLL shared task, systems are also evaluated on the out-of-domain dataset. Statistical models are typically less accurate when they are applied to out-of-domain data. Consequently, the predicted syntax for the out-of-domain test set is of lower quality, which negatively affects the quality of GCN embeddings. However, our model works surprisingly well on out-of-domain data (Table 5), substantially outperforming all the previous syntax-aware models. This suggests that our model is fairly robust to mistakes in syntax. As expected though, our model does not outperform the syntax-agnostic model of Marcheggiani et al. (2017).

## 6 Related Work

Perhaps the earliest methods modeling syntax-semantics interface with RNNs are due to (Henderson et al., 2008; Titov et al., 2009; Gesmundo et al., 2009), they used shift-reduce parsers for joint SRL and syntactic parsing, and relied on RNNs to model statistical dependencies across syntactic and semantic parsing actions. A more

<sup>9</sup>To compare to previous work, we report combined scores which also include predicate disambiguation. As we use disambiguators from previous work (see Section 5.1), actual gains in argument identification and labeling are even larger.

modern (e.g., based on LSTMs) and effective reincarnation of this line of research has been proposed in Swayamdipta et al. (2016). Other recent work which considered incorporation of syntactic information in neural SRL models include: FitzGerald et al. (2015) who use standard syntactic features within an MLP calculating potentials of a CRF model; Roth and Lapata (2016) who enriched standard features for SRL with LSTM representations of syntactic paths between arguments and predicates; Lei et al. (2015) who relied on low-rank tensor factorizations for modeling syntax. Also Foland and Martin (2015) used (non-graph) convolutional networks and provided syntactic features as input. A very different line of research, but with similar goals to ours (i.e. integrating syntax with minimal feature engineering), used tree kernels (Moschitti et al., 2008).

Beyond SRL, there have been many proposals on how to incorporate syntactic information in RNN models, for example, in the context of neural machine translation (Eriguchi et al., 2017; Sennrich and Haddow, 2016). One of the most popular and attractive approaches is to use tree-structured recursive neural networks (Socher et al., 2013; Le and Zuidema, 2014; Dyer et al., 2015), including stacking them on top of a sequential BiLSTM (Miwa and Bansal, 2016). An approach of Mou et al. (2015) to sentiment analysis and question classification, introduced even before GCNs became popular in the machine learning community, is related to graph convolution. However, it is inherently single-layer and tree-specific, uses bottom-up computations, does not share parameters across syntactic functions and does not use gates. Gates have been previously used in GCNs (Li et al., 2016) but between GCN layers rather than for individual edges.

Previous approaches to integrating syntactic information in neural models are mainly designed to induce representations of sentences or syntactic constituents. In contrast, the approach we presented incorporates syntactic information at word level. This may be attractive from the engineering perspective, as it can be used, as we have shown, instead or along with RNN models.

## 7 Conclusions and Future Work

We demonstrated how GCNs can be used to incorporate syntactic information in neural models and specifically to construct a syntax-aware SRL



model, resulting in state-of-the-art results for Chinese and English. There are relatively straightforward steps which can further improve the SRL results. **For example, we relied on labeling arguments independently, whereas using a joint model is likely to significantly improve the performance.**

More generally, given simplicity of GCNs and their applicability to general graph structures (not necessarily trees), we believe that there are many NLP tasks where GCNs can be used to incorporate linguistic structures (e.g., syntactic and semantic representations of sentences and discourse parses or co-reference graphs for documents).

## Acknowledgements

We would thank Anton Frolov, Michael Schlichtkrull, Thomas Kipf, Michael Roth, Max Welling, Yi Zhang, and Wilker Aziz for their suggestions and comments. The project was supported by the European Research Council (ERC StG BroadSem 678254), the Dutch National Science Foundation (NWO VIDI 639.022.518) and an Amazon Web Services (AWS) grant.

## References

- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. A neural probabilistic language model. *Journal of Machine Learning Research* 3:1137–1155.
- Anders Björkelund, Bernd Bohnet, Love Hafdell, and Pierre Nugues. 2010. A high-performance syntactic and semantic dependency parser. In *Proceedings of COLING: Demonstrations*.
- Anders Björkelund, Love Hafdell, and Pierre Nugues. 2009. Multilingual semantic role labeling. In *Proceedings of CoNLL*.
- Jason P. C. Chiu and Eric Nichols. 2016. Named entity recognition with bidirectional LSTM-CNNs. *TACL* 4:357–370.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of EMNLP*.
- Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. 2016. Language modeling with gated convolutional networks. *arXiv preprint arXiv:1612.08083*.
- David K. Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P. Adams. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *Proceedings of NIPS*.
- 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 ACL*.
- Jeffrey L. Elman. 1990. Finding structure in time. *Cognitive Science* 14(2):179–211.
- Akiko Eriguchi, Yoshimasa Tsuruoka, and Kyunghyun Cho. 2017. Learning to parse and translate improves neural machine translation. *arXiv preprint arXiv:1702.03525*.
- Nicholas FitzGerald, Oscar Täckström, Kuzman Ganchev, and Dipanjan Das. 2015. Semantic role labeling with neural network factors. In *Proceedings of EMNLP*.
- William Foland and James Martin. 2015. Dependency-based semantic role labeling using convolutional neural networks. In *Proceedings of the Fourth Joint Conference on Lexical and Computational Semantics (\*SEM)*.
- Andrea Gesmundo, James Henderson, Paola Merlo, and Ivan Titov. 2009. Latent variable model of synchronous syntactic-semantic parsing for multiple languages. In *Proceedings of CoNLL*.
- Daniel Gildea and Daniel Jurafsky. 2002. Automatic labeling of semantic roles. *Computational linguistics* 28(3):245–288.
- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. 2017. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*.
- Alex Graves. 2008. *Supervised sequence labelling with recurrent neural networks*. Ph.D. thesis, München, Techn. Univ., Diss., 2008.
- James Henderson, Paola Merlo, Gabriele Musillo, and Ivan Titov. 2008. A latent variable model of synchronous parsing for syntactic and semantic dependencies. In *Proceedings of CoNLL*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation* 9(8):1735–1780.
- Richard Johansson and Pierre Nugues. 2008. The effect of syntactic representation on semantic role labeling. In *Proceedings of COLING*.
- Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay Pande, and Patrick Riley. 2016. Molecular graph convolutions: moving beyond fingerprints. *Journal of computer-aided molecular design* 30(8):595–608.
- Diederik Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proceedings of ICLR*.

- Eliyahu Kiperwasser and Yoav Goldberg. 2016. Simple and accurate dependency parsing using bidirectional LSTM feature representations. *TACL* 4:313–327.
- Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *Proceedings of ICLR*.
- Phong Le and Willem Zuidema. 2014. The inside-outside recursive neural network model for dependency parsing. In *Proceedings of EMNLP*.
- Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. 2001. Gradient-based learning applied to document recognition. In *Proceedings of Intelligent Signal Processing*.
- Tao Lei, Yuan Zhang, Lluís Màrquez, Alessandro Moschitti, and Regina Barzilay. 2015. High-order low-rank tensors for semantic role labeling. In *Proceedings of NAACL*.
- Beth Levin. 1993. *English verb classes and alternations: A preliminary investigation*. University of Chicago press.
- Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. 2016. Gated graph sequence neural networks. In *Proceedings of ICLR*.
- Wang Ling, Chris Dyer, Alan W Black, and Isabel Trancoso. 2015. Two/too simple adaptations of word2vec for syntax problems. In *Proceedings of NAACL*.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *TACL* 4:521–535.
- Diego Marcheggiani, Anton Frolov, and Ivan Titov. 2017. A simple and accurate syntax-agnostic neural model for dependency-based semantic role labeling. *arXiv preprint arXiv:1701.02593*.
- Tomas Mikolov, Martin Karafiát, Lukás Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Proceedings of INTERSPEECH*.
- Makoto Miwa and Mohit Bansal. 2016. End-to-end relation extraction using lstms on sequences and tree structures. In *Proceedings of ACL*.
- Alessandro Moschitti, Daniele Pighin, and Roberto Basili. 2008. Tree kernels for semantic role labeling. *Computational Linguistics* 34(2):193–224.
- Lili Mou, Hao Peng, Ge Li, Yan Xu, Lu Zhang, and Zhi Jin. 2015. Discriminative neural sentence modeling by tree-based convolution. In *Proceedings of EMNLP*.
- Sameer Pradhan, Kadri Hacioglu, Wayne H. Ward, James H. Martin, and Daniel Jurafsky. 2005. Semantic role chunking combining complementary syntactic views. In *Proceedings of CoNLL*.
- Vasin Punyakanok, Dan Roth, and Wen-tau Yih. 2008. The importance of syntactic parsing and inference in semantic role labeling. *Computational Linguistics* 34(2):257–287.
- Michael Roth and Mirella Lapata. 2016. Neural semantic role labeling with dependency path embeddings. In *Proceedings of ACL*.
- Rico Sennrich and Barry Haddow. 2016. Linguistic input features improve neural machine translation. In *Proceedings of WMT*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of EMNLP*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of NIPS*.
- Swabha Swayamdipta, Miguel Ballesteros, Chris Dyer, and Noah A. Smith. 2016. Greedy, joint syntactic-semantic parsing with stack LSTMs. In *Proceedings of CoNLL*.
- Cynthia A. Thompson, Roger Levy, and Christopher D. Manning. 2003. A generative model for semantic role labeling. In *Proceedings of ECML*.
- Ivan Titov, James Henderson, Paola Merlo, and Gabriele Musillo. 2009. Online projectivisation for synchronous parsing of semantic and syntactic dependencies. In *Proceedings of IJCAI*.
- Aäron van den Oord, Nal Kalchbrenner, Lasse Espeholt, Koray Kavukcuoglu, Oriol Vinyals, and Alex Graves. 2016. Conditional image generation with PixelCNN decoders. In *Proceedings of NIPS*.
- Hai Zhao, Wenliang Chen, Jun’ichi Kazama, Kiyotaka Uchimoto, and Kentaro Torisawa. 2009. Multilingual dependency learning: Exploiting rich features for tagging syntactic and semantic dependencies. In *Proceedings of CoNLL*.
- Jie Zhou and Wei Xu. 2015. End-to-end learning of semantic role labeling using recurrent neural networks. In *Proceedings of ACL*.

## A Hyperparameter values

| Semantic role labeler                |     |
|--------------------------------------|-----|
| $d_w$ (word embeddings EN)           | 100 |
| $d_w$ (word embeddings CH)           | 128 |
| $d_{pos}$ (POS embeddings)           | 16  |
| $d_l$ (lemma embeddings)             | 100 |
| $d_h$ (LSTM hidden states)           | 512 |
| $d_r$ (role representation)          | 128 |
| $d'_l$ (output lemma representation) | 128 |
| $J$ (BiLSTM depth)                   | 3   |
| $K$ (GCN depth)                      | 1   |
| $\beta$ (edge dropout)               | .3  |
| learning rate                        | .01 |

Table 6: Hyperparameter values.