

# A Joint Many-Task Model: Growing a Neural Network for Multiple NLP Tasks

Kazuma Hashimoto\*, Caiming Xiong<sup>†</sup>, Yoshimasa Tsuruoka, and Richard Socher

The University of Tokyo

{hassy, tsuruoka}@logos.t.u-tokyo.ac.jp

Salesforce Research

{cxiong, rsocher}@salesforce.com

## Abstract

Transfer and multi-task learning have traditionally focused on either a single source-target pair or very few, similar tasks. Ideally, the linguistic levels of morphology, syntax and semantics would benefit each other by being trained in a single model. We introduce a joint many-task model together with a strategy for successively growing its depth to solve increasingly complex tasks. Higher layers include shortcut connections to lower-level task predictions to reflect linguistic hierarchies. We use a simple regularization term to allow for optimizing all model weights to improve one task’s loss without exhibiting catastrophic interference of the other tasks. Our single **end-to-end model** obtains state-of-the-art or competitive results on five different tasks from tagging, parsing, relatedness, and entailment tasks.

## 1 Introduction

The potential for leveraging multiple levels of representation has been demonstrated in various ways in the field of Natural Language Processing (NLP). For example, Part-Of-Speech (POS) tags are used for syntactic parsers. **The parsers are used to improve higher-level tasks, such as natural language inference (Chen et al., 2016) and machine translation (Eriguchi et al., 2016).** These systems are often pipelines and not trained end-to-end.

Deep NLP models have yet shown benefits from predicting many increasingly complex tasks each at a successively deeper layer. Existing models often ignore linguistic hierarchies by predicting

\* Work was done while the first author was an intern at Salesforce Research.

<sup>†</sup> Corresponding author.

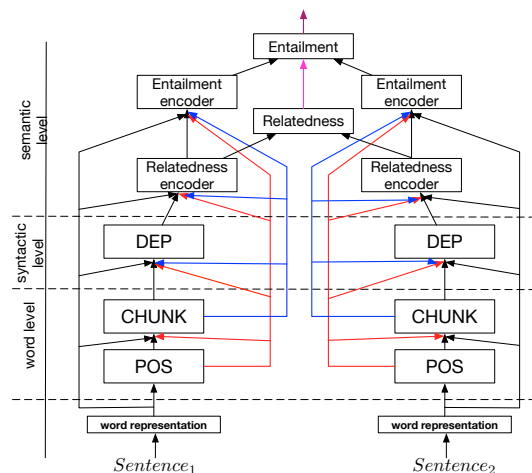


Figure 1: Overview of the joint many-task model predicting different linguistic outputs at successively deeper layers.

different tasks either entirely separately or at the same depth (Collobert et al., 2011).

We introduce a Joint **Many-Task (JMT) model**, outlined in Figure 1, which predicts increasingly complex NLP tasks at successively deeper layers. Unlike traditional pipeline systems, our single JMT model can be trained end-to-end for **POS tagging, chunking, dependency parsing, semantic relatedness, and textual entailment, by considering linguistic hierarchies.** We propose an adaptive training and regularization strategy to grow this model in its depth. With the help of this strategy we avoid catastrophic interference between the tasks. Our model is motivated by Søgaard and Goldberg (2016) who showed that predicting two different tasks is more accurate when performed in different layers than in the same layer (Collobert et al., 2011). Experimental results show that our single model achieves competitive results for all of the five different tasks, demonstrating that us-

ing linguistic hierarchies is more important than handling different tasks in the same layer.

## 2 The Joint Many-Task Model

This section describes the inference procedure of our model, beginning at the lowest level and working our way to higher layers and more complex tasks; our model handles the five different tasks in the order of POS tagging, chunking, dependency parsing, semantic relatedness, and textual entailment, by considering linguistic hierarchies. The POS tags are used for chunking, and the chunking tags are used for dependency parsing (Attardi and Dell’Orletta, 2008). Tai et al. (2015) have shown that dependencies improve the relatedness task. The relatedness and entailment tasks are closely related to each other. If the semantic relatedness between two sentences is very low, they are unlikely to entail each other. Based on this observation, we make use of the information from the relatedness task for improving the entailment task.

### 2.1 Word Representations

For each word  $w_t$  in the input sentence  $s$  of length  $L$ , we use two types of embeddings.

**Word embeddings:** We use Skip-gram (Mikolov et al., 2013) to train word embeddings.

**Character embeddings:** Character  $n$ -gram embeddings are trained by the same Skip-gram objective. We construct the character  $n$ -gram vocabulary in the training data and assign an embedding for each entry. The final character embedding is the average of the *unique* character  $n$ -gram embeddings of  $w_t$ . For example, the character  $n$ -grams ( $n = 1, 2, 3$ ) of the word “Cat” are {C, a, t, #B#C, Ca, at, t#E#, #B#Ca, Cat, at#E#}, where “#B#” and “#E#” represent the beginning and the end of each word, respectively. Using the character embeddings efficiently provides morphological features. Each word is subsequently represented as  $x_t$ , the concatenation of its corresponding word and character embeddings shared across the tasks.<sup>1</sup>

### 2.2 Word-Level Task: POS Tagging

The first layer of the model is a bi-directional LSTM (Graves and Schmidhuber, 2005; Hochreiter and Schmidhuber, 1997) whose hidden states

are used to predict POS tags. We use the following Long Short-Term Memory (LSTM) units for the forward direction:

$$\begin{aligned} i_t &= \sigma(W_i g_t + b_i), \quad f_t = \sigma(W_f g_t + b_f), \\ o_t &= \sigma(W_o g_t + b_o), \quad u_t = \tanh(W_u g_t + b_u), \\ c_t &= i_t \odot u_t + f_t \odot c_{t-1}, \quad h_t = o_t \odot \tanh(c_t), \end{aligned} \quad (1)$$

where we define the input  $g_t$  as  $g_t = [\vec{h}_{t-1}; x_t]$ , i.e. the concatenation of the previous hidden state and the word representation of  $w_t$ . The backward pass is expanded in the same way, but a different set of weights are used.

For predicting the POS tag of  $w_t$ , we use the concatenation of the forward and backward states in a one-layer bi-LSTM layer corresponding to the  $t$ -th word:  $h_t = [\vec{h}_t; \overleftarrow{h}_t]$ . Then each  $h_t$  ( $1 \leq t \leq L$ ) is fed into a standard softmax classifier with a single ReLU layer which outputs the probability vector  $y^{(1)}$  for each of the POS tags.

### 2.3 Word-Level Task: Chunking

Chunking is also a word-level classification task which assigns a chunking tag (B-NP, I-VP, etc.) for each word. The tag specifies the region of major phrases (e.g., noun phrases) in the sentence.

Chunking is performed in the second bi-LSTM layer on top of the POS layer. When stacking the bi-LSTM layers, we use Eq. (1) with input  $g_t^{(2)} = [h_{t-1}^{(2)}; h_t^{(1)}; x_t; y_t^{(pos)}]$ , where  $h_t^{(1)}$  is the hidden state of the first (POS) layer. We define the weighted label embedding  $y_t^{(pos)}$  as follows:

$$y_t^{(pos)} = \sum_{j=1}^C p(y_t^{(1)} = j | h_t^{(1)}) \ell(j), \quad (2)$$

where  $C$  is the number of the POS tags,  $p(y_t^{(1)} = j | h_t^{(1)})$  is the probability value that the  $j$ -th POS tag is assigned to  $w_t$ , and  $\ell(j)$  is the corresponding label embedding. The probability values are predicted by the POS layer, and thus no gold POS tags are needed. This output embedding is similar to the  $K$ -best POS tag feature which has been shown to be effective in syntactic tasks (Andor et al., 2016; Alberti et al., 2015). For predicting the chunking tags, we employ the same strategy as POS tagging by using the concatenated bi-directional hidden states  $h_t^{(2)} = [\vec{h}_t^{(2)}; \overleftarrow{h}_t^{(2)}]$  in the chunking layer. We also use a single ReLU hidden layer before the softmax classifier.

<sup>1</sup>Bojanowski et al. (2017) previously proposed to train the character  $n$ -gram embeddings by the Skip-gram objective.

## 2.4 Syntactic Task: Dependency Parsing

Dependency parsing identifies syntactic relations (such as an adjective modifying a noun) between word pairs in a sentence. We use the third bi-LSTM layer to classify relations between all pairs of words. The input vector for the LSTM includes hidden states, word representations, and the label embeddings for the two previous tasks:  $g_t^{(3)} = [h_{t-1}^{(3)}; h_t^{(2)}; x_t; (y_t^{(pos)} + y_t^{(chk)})]$ , where we computed the chunking vector in a similar fashion as the POS vector in Eq. (2).

We predict the parent node (*head*) for each word. Then a dependency label is predicted for each child-parent pair. This approach is related to Dozat and Manning (2017) and Zhang et al. (2017), where the main difference is that our model works on a multi-task framework. To predict the parent node of  $w_t$ , we define a matching function between  $w_t$  and the candidates of the parent node as  $m(t, j) = h_t^{(3)} \cdot (W_d h_j^{(3)})$ , where  $W_d$  is a parameter matrix. For the root, we define  $h_{L+1}^{(3)} = r$  as a parameterized vector. To compute the probability that  $w_j$  (or the root node) is the parent of  $w_t$ , the scores are normalized:

$$p(j|h_t^{(3)}) = \frac{\exp(m(t, j))}{\sum_{k=1, k \neq t}^{L+1} \exp(m(t, k))}. \quad (3)$$

The dependency labels are predicted using  $[h_t^{(3)}; h_j^{(3)}]$  as input to a softmax classifier with a single ReLU layer. We greedily select the parent node and the dependency label for each word. When the parsing result is not a well-formed tree, we apply the first-order Eisner’s algorithm (Eisner, 1996) to obtain a well-formed tree from it.

## 2.5 Semantic Task: Semantic relatedness

The next two tasks model the semantic relationships between two input sentences. The first task measures the semantic relatedness between two sentences. The output is a real-valued relatedness score for the input sentence pair. The second task is textual entailment, which requires one to determine whether a premise sentence entails a hypothesis sentence. There are typically three classes: entailment, contradiction, and neutral. We use the fourth and fifth bi-LSTM layer for the relatedness and entailment task, respectively.

Now it is required to obtain the sentence-level representation rather than the word-level representation  $h_t^{(4)}$  used in the first three tasks. We compute the sentence-level representation  $h_s^{(4)}$  as the

element-wise maximum values across all of the word-level representations in the fourth layer:

$$h_s^{(4)} = \max(h_1^{(4)}, h_2^{(4)}, \dots, h_L^{(4)}). \quad (4)$$

This max-pooling technique has proven effective in text classification tasks (Lai et al., 2015).

To model the semantic relatedness between  $s$  and  $s'$ , we follow Tai et al. (2015). The feature vector for representing the semantic relatedness is computed as follows:

$$d_1(s, s') = \left[ |h_s^{(4)} - h_{s'}^{(4)}|; h_s^{(4)} \odot h_{s'}^{(4)} \right], \quad (5)$$

where  $|h_s^{(4)} - h_{s'}^{(4)}|$  is the absolute values of the element-wise subtraction, and  $h_s^{(4)} \odot h_{s'}^{(4)}$  is the element-wise multiplication. Then  $d_1(s, s')$  is fed into a softmax classifier with a single Maxout hidden layer (Goodfellow et al., 2013) to output a relatedness score (from 1 to 5 in our case).

## 2.6 Semantic Task: Textual entailment

For entailment classification, we also use the max-pooling technique as in the semantic relatedness task. To classify the premise-hypothesis pair  $(s, s')$  into one of the three classes, we compute the feature vector  $d_2(s, s')$  as in Eq. (5) except that we do not use the absolute values of the element-wise subtraction, because we need to identify which is the premise (or hypothesis). Then  $d_2(s, s')$  is fed into a softmax classifier.

To use the output from the relatedness layer directly, we use the label embeddings for the relatedness task. More concretely, we compute the class label embeddings for the semantic relatedness task similar to Eq. (2). The final feature vectors that are concatenated and fed into the entailment classifier are the weighted relatedness label embedding and the feature vector  $d_2(s, s')$ . We use three Maxout hidden layers before the classifier.

## 3 Training the JMT Model

The model is trained jointly over all datasets. During each epoch, the optimization iterates over each full training dataset in the same order as the corresponding tasks described in the modeling section.

### 3.1 Pre-Training Word Representations

We pre-train word embeddings using the Skip-gram model with negative sampling (Mikolov

et al., 2013). We also pre-train the character  $n$ -gram embeddings using Skip-gram.<sup>2</sup> The only difference is that each input word embedding is replaced with its corresponding average character  $n$ -gram embedding described in Section 2.1. These embeddings are fine-tuned during the model training. We denote the embedding parameters as  $\theta_e$ .

### 3.2 Training the POS Layer

Let  $\theta_{\text{POS}} = (W_{\text{POS}}, b_{\text{POS}}, \theta_e)$  denote the set of model parameters associated with the POS layer, where  $W_{\text{POS}}$  is the set of the weight matrices in the first bi-LSTM and the classifier, and  $b_{\text{POS}}$  is the set of the bias vectors. The objective function to optimize  $\theta_{\text{POS}}$  is defined as follows:

$$J_1(\theta_{\text{POS}}) = - \sum_s \sum_t \log p(y_t^{(1)} = \alpha | h_t^{(1)}) + \lambda \|W_{\text{POS}}\|^2 + \delta \|\theta_e - \theta'_e\|^2, \quad (6)$$

where  $p(y_t^{(1)} = \alpha | h_t^{(1)})$  is the probability value that the correct label  $\alpha$  is assigned to  $w_t$  in the sentence  $s$ ,  $\lambda \|W_{\text{POS}}\|^2$  is the L2-norm regularization term, and  $\lambda$  is a hyperparameter.

We call the second regularization term  $\delta \|\theta_e - \theta'_e\|^2$  a *successive* regularization term. The successive regularization is based on the idea that we do not want the model to forget the information learned for the other tasks. In the case of POS tagging, the regularization is applied to  $\theta_e$ , and  $\theta'_e$  is the embedding parameter after training the final task in the top-most layer at the previous training epoch.  $\delta$  is a hyperparameter.

### 3.3 Training the Chunking Layer

The objective function is defined as follows:

$$J_2(\theta_{\text{chk}}) = - \sum_s \sum_t \log p(y_t^{(2)} = \alpha | h_t^{(2)}) + \lambda \|W_{\text{chk}}\|^2 + \delta \|\theta_{\text{POS}} - \theta'_{\text{POS}}\|^2, \quad (7)$$

which is similar to that of POS tagging, and  $\theta_{\text{chk}}$  is  $(W_{\text{chk}}, b_{\text{chk}}, E_{\text{POS}}, \theta_e)$ , where  $W_{\text{chk}}$  and  $b_{\text{chk}}$  are the weight and bias parameters including those in  $\theta_{\text{POS}}$ , and  $E_{\text{POS}}$  is the set of the POS label embeddings.  $\theta'_{\text{POS}}$  is the one after training the POS layer at the current training epoch.

<sup>2</sup>The training code and the pre-trained embeddings are available at <https://github.com/hassyGo/charNgram2vec>.

### 3.4 Training the Dependency Layer

The objective function is defined as follows:

$$J_3(\theta_{\text{dep}}) = - \sum_s \sum_t \log p(\alpha | h_t^{(3)}) p(\beta | h_t^{(3)}, h_\alpha^{(3)}) + \lambda (\|W_{\text{dep}}\|^2 + \|W_d\|^2) + \delta \|\theta_{\text{chk}} - \theta'_{\text{chk}}\|^2, \quad (8)$$

where  $p(\alpha | h_t^{(3)})$  is the probability value assigned to the correct parent node  $\alpha$  for  $w_t$ , and  $p(\beta | h_t^{(3)}, h_\alpha^{(3)})$  is the probability value assigned to the correct dependency label  $\beta$  for the child-parent pair  $(w_t, \alpha)$ .  $\theta_{\text{dep}}$  is defined as  $(W_{\text{dep}}, b_{\text{dep}}, W_d, r, E_{\text{POS}}, E_{\text{chk}}, \theta_e)$ , where  $W_{\text{dep}}$  and  $b_{\text{dep}}$  are the weight and bias parameters including those in  $\theta_{\text{chk}}$ , and  $E_{\text{chk}}$  is the set of the chunking label embeddings.

### 3.5 Training the Relatedness Layer

Following Tai et al. (2015), the objective function is defined as follows:

$$J_4(\theta_{\text{rel}}) = \sum_{(s,s')} \text{KL} \left( \hat{p}(s, s') \parallel p(h_s^{(4)}, h_{s'}^{(4)}) \right) + \lambda \|W_{\text{rel}}\|^2 + \delta \|\theta_{\text{dep}} - \theta'_{\text{dep}}\|^2, \quad (9)$$

where  $\hat{p}(s, s')$  is the gold distribution over the defined relatedness scores,  $p(h_s^{(4)}, h_{s'}^{(4)})$  is the predicted distribution given the sentence representations, and  $\text{KL} \left( \hat{p}(s, s') \parallel p(h_s^{(4)}, h_{s'}^{(4)}) \right)$  is the KL-divergence between the two distributions.  $\theta_{\text{rel}}$  is defined as  $(W_{\text{rel}}, b_{\text{rel}}, E_{\text{POS}}, E_{\text{chk}}, \theta_e)$ .

### 3.6 Training the Entailment Layer

The objective function is defined as follows:

$$J_5(\theta_{\text{ent}}) = - \sum_{(s,s')} \log p(y_{(s,s')}^{(5)} = \alpha | h_s^{(5)}, h_{s'}^{(5)}) + \lambda \|W_{\text{ent}}\|^2 + \delta \|\theta_{\text{rel}} - \theta'_{\text{rel}}\|^2, \quad (10)$$

where  $p(y_{(s,s')}^{(5)} = \alpha | h_s^{(5)}, h_{s'}^{(5)})$  is the probability value that the correct label  $\alpha$  is assigned to the premise-hypothesis pair  $(s, s')$ .  $\theta_{\text{ent}}$  is defined as  $(W_{\text{ent}}, b_{\text{ent}}, E_{\text{POS}}, E_{\text{chk}}, E_{\text{rel}}, \theta_e)$ , where  $E_{\text{rel}}$  is the set of the relatedness label embeddings.

## 4 Related Work

Many deep learning approaches have proven to be effective in a variety of NLP tasks and are becoming more and more complex. They are typically



designed to handle single tasks, or some of them are designed as general-purpose models (Kumar et al., 2016; Sutskever et al., 2014) but applied to different tasks independently.

For handling multiple NLP tasks, multi-task learning models with deep neural networks have been proposed (Collobert et al., 2011; Luong et al., 2016), and more recently Søgaard and Goldberg (2016) have suggested that using different layers for different tasks is more effective than using the same layer in jointly learning closely-related tasks, such as POS tagging and chunking. However, the number of tasks was limited or they have very similar task settings like word-level tagging, and it was not clear how lower-level tasks could be also improved by combining higher-level tasks.

More related to our work, Godwin et al. (2016) also followed Søgaard and Goldberg (2016) to jointly learn POS tagging, chunking, and language modeling, and Zhang and Weiss (2016) have shown that it is effective to jointly learn POS tagging and dependency parsing by sharing internal representations. In the field of relation extraction, Miwa and Bansal (2016) proposed a joint learning model for entity detection and relation extraction. All of them suggest the importance of multi-task learning, and we investigate the potential of handling different types of NLP tasks rather than closely-related ones in a single hierarchical deep model.

In the field of computer vision, some transfer and multi-task learning approaches have also been proposed (Li and Hoiem, 2016; Misra et al., 2016). For example, Misra et al. (2016) proposed a multi-task learning model to handle different tasks. However, they assume that each data sample has annotations for the different tasks, and do not explicitly consider task hierarchies.

Recently, Rusu et al. (2016) have proposed a progressive neural network model to handle multiple reinforcement learning tasks, such as Atari games. Like our JMT model, their model is also successively trained according to different tasks using different layers called columns in their paper. In their model, once the first task is completed, the model parameters for the first task are fixed, and then the second task is handled with new model parameters. Therefore, accuracy of the previously trained tasks is never improved. In NLP tasks, multi-task learning has the potential to improve not only higher-level tasks, but also lower-

level tasks. Rather than fixing the pre-trained model parameters, our successive regularization allows our model to continuously train the lower-level tasks without significant accuracy drops.

## 5 Experimental Settings

### 5.1 Datasets

**POS tagging:** To train the POS tagging layer, we used the Wall Street Journal (WSJ) portion of Penn Treebank, and followed the standard split for the training (Section 0-18), development (Section 19-21), and test (Section 22-24) sets. The evaluation metric is the word-level accuracy.

**Chunking:** For chunking, we also used the WSJ corpus, and followed the standard split for the training (Section 15-18) and test (Section 20) sets as in the CoNLL 2000 shared task. We used Section 19 as the development set and employed the IOBES tagging scheme. The evaluation metric is the F1 score defined in the shared task.

**Dependency parsing:** We also used the WSJ corpus for dependency parsing, and followed the standard split for the training (Section 2-21), development (Section 22), and test (Section 23) sets. We obtained Stanford style dependencies using the version 3.3.0 of the Stanford converter. The evaluation metrics are the Unlabeled Attachment Score (UAS) and the Labeled Attachment Score (LAS), and punctuations are excluded for the evaluation.

**Semantic relatedness:** For the semantic relatedness task, we used the SICK dataset (Marelli et al., 2014), and followed the standard split for the training, development, and test sets. The evaluation metric is the Mean Squared Error (MSE) between the gold and predicted scores.

**Textual entailment:** For textual entailment, we also used the SICK dataset and exactly the same data split as the semantic relatedness dataset. The evaluation metric is the accuracy.

### 5.2 Training Details

We set the dimensionality of the embeddings and the hidden states in the bi-LSTMs to 100. At each training epoch, we trained our model in the order of POS tagging, chunking, dependency parsing, semantic relatedness, and textual entailment. We used mini-batch stochastic gradient descent and empirically found it effective to use a gradient clipping method with growing clipping values for the different tasks; concretely, we employed the simple function:  $\min(3.0, depth)$ , where  $depth$  is

the number of bi-LSTM layers involved in each task, and 3.0 is the maximum value. We applied our successive regularization to our model, along with L2-norm regularization and dropout (Srivastava et al., 2014). More details are summarized in the supplemental material.

## 6 Results and Discussion

Table 1 shows our results on the test sets of the five tasks.<sup>3</sup> The column “Single” shows the results of handling each task separately using single-layer bi-LSTMs, and the column “JMT<sub>all</sub>” shows the results of our JMT model. The single task settings only use the annotations of their own tasks. For example, when handling dependency parsing as a single task, the POS and chunking tags are *not* used. We can see that all results of the five tasks are improved in our JMT model, which shows that our JMT model can handle the five different tasks in a single model. Our JMT model allows us to access arbitrary information learned from the different tasks. If we want to use the model just as a POS tagger, we can use only first bi-LSTM layer.

Table 1 also shows the results of five subsets of the different tasks. For example, in the case of “JMT<sub>ABC</sub>”, only the first three layers of the bi-LSTMs are used to handle the three tasks. In the case of “JMT<sub>DE</sub>”, only the top two layers are used as a two-layer bi-LSTM by omitting all information from the first three layers. The results of the closely-related tasks (“AB”, “ABC”, and “DE”) show that our JMT model improves both of the high-level and low-level tasks. The results of “JMT<sub>CD</sub>” and “JMT<sub>CE</sub>” show that the parsing task can be improved by the semantic tasks.

It should be noted that in our analysis on the greedy parsing results of the “JMT<sub>ABC</sub>” setting, we have found that more than 95% are well-formed dependency trees on the development set. In the 1,700 sentences of the development data, 11 results have multiple root nodes, 11 results have no root nodes, and 61 results have cycles. These 83 parsing results are converted into well-formed trees by Eisner’s algorithm, and the accuracy does not significantly change (UAS: 94.52%→94.53%, LAS: 92.61%→92.62%).

<sup>3</sup>In chunking evaluation, we only show the results of “Single” and “JMT<sub>AB</sub>” because the sentences for chunking evaluation overlap the training data for dependency parsing.

### 6.1 Comparison with Published Results

**POS tagging** Table 2 shows the results of POS tagging, and our JMT model achieves the score close to the state-of-the-art results. The best result to date has been achieved by Ling et al. (2015), which uses character-based LSTMs. Incorporating the character-based encoders into our JMT model would be an interesting direction, but we have shown that the simple pre-trained character  $n$ -gram embeddings lead to the promising result.

**Chunking** Table 3 shows the results of chunking, and our JMT model achieves the state-of-the-art result. Søgaard and Goldberg (2016) proposed to jointly learn POS tagging and chunking in different layers, but they only showed improvement for chunking. By contrast, our results show that the low-level tasks are also improved.

**Dependency parsing** Table 4 shows the results of dependency parsing by using only the WSJ corpus in terms of the dependency annotations.<sup>4</sup> It is notable that our simple greedy dependency parser outperforms the model in Andor et al. (2016) which is based on beam search with global information. The result suggests that the bi-LSTMs efficiently capture global information necessary for dependency parsing. Moreover, our single task result already achieves high accuracy without the POS and chunking information. The best result to date has been achieved by the model proposed in Dozat and Manning (2017), which uses higher dimensional representations than ours and proposes a more sophisticated attention mechanism called *biaffine attention*. It should be promising to incorporate their attention mechanism into our parsing component.

**Semantic relatedness** Table 5 shows the results of the semantic relatedness task, and our JMT model achieves the state-of-the-art result. The result of “JMT<sub>DE</sub>” is already better than the previous state-of-the-art results. Both of Zhou et al. (2016) and Tai et al. (2015) explicitly used syntactic trees, and Zhou et al. (2016) relied on attention mechanisms. However, our method uses the simple max-pooling strategy, which suggests that it is worth

<sup>4</sup>Choe and Charniak (2016) employed a tri-training method to expand the training data with 400,000 trees in addition to the WSJ data, and they reported 95.9 UAS and 94.1 LAS by converting their constituency trees into dependency trees. Kuncoro et al. (2017) also reported high accuracy (95.8 UAS and 94.6 LAS) by using a converter.

		Single	JMT <sub>all</sub>	JMT <sub>AB</sub>	JMT <sub>ABC</sub>	JMT <sub>DE</sub>	JMT <sub>CD</sub>	JMT <sub>CE</sub>
A ↑	POS	97.45	97.55	97.52	97.54	n/a	n/a	n/a
B ↑	Chunking	95.02	n/a	95.77	n/a	n/a	n/a	n/a
C ↑	Dependency UAS	93.35	94.67	n/a	94.71	n/a	93.53	93.57
	Dependency LAS	91.42	92.90	n/a	92.92	n/a	91.62	91.69
D ↓	Relatedness	0.247	0.233	n/a	n/a	0.238	0.251	n/a
E ↑	Entailment	81.8	86.2	n/a	n/a	86.8	n/a	82.4

Table 1: Test set results for the five tasks. In the relatedness task, the lower scores are better.

Method	Acc. ↑
JMT <sub>all</sub>	97.55
Ling et al. (2015)	<b>97.78</b>
Kumar et al. (2016)	97.56
Ma and Hovy (2016)	97.55
Søgaard (2011)	97.50
Collobert et al. (2011)	97.29
Tsuruoka et al. (2011)	97.28
Toutanova et al. (2003)	97.27

Table 2: POS tagging results.

Method	F1 ↑
JMT <sub>AB</sub>	<b>95.77</b>
Single	95.02
Søgaard and Goldberg (2016)	95.56
Suzuki and Isozaki (2008)	95.15
Collobert et al. (2011)	94.32
Kudo and Matsumoto (2001)	93.91
Tsuruoka et al. (2011)	93.81

Table 3: Chunking results.

Method	UAS ↑	LAS ↑
JMT <sub>all</sub>	94.67	92.90
Single	93.35	91.42
Dozat and Manning (2017)	<b>95.74</b>	<b>94.08</b>
Andor et al. (2016)	94.61	92.79
Alberti et al. (2015)	94.23	92.36
Zhang et al. (2017)	94.10	91.90
Weiss et al. (2015)	93.99	92.05
Dyer et al. (2015)	93.10	90.90
Bohnet (2010)	92.88	90.71

Table 4: Dependency results.

Method	MSE ↓
JMT <sub>all</sub>	<b>0.233</b>
JMT <sub>DE</sub>	0.238
Zhou et al. (2016)	0.243
Tai et al. (2015)	0.253

Table 5: Semantic relatedness results.

Method	Acc. ↑
JMT <sub>all</sub>	86.2
JMT <sub>DE</sub>	<b>86.8</b>
Yin et al. (2016)	86.2
Lai and Hockenmaier (2014)	84.6

Table 6: Textual entailment results.

	JMT <sub>all</sub>	w/o SC	w/o LE	w/o SC&LE
POS	97.88	97.79	97.85	97.87
Chunking	97.59	97.08	97.40	97.33
Dependency UAS	94.51	94.52	94.09	94.04
Dependency LAS	92.60	92.62	92.14	92.03
Relatedness	0.236	0.698	0.261	0.765
Entailment	84.6	75.0	81.6	71.2

Table 7: Effectiveness of the Shortcut Connections (SC) and the Label Embeddings (LE).

	JMT <sub>ABC</sub>	w/o SC&LE	All-3
POS	97.90	97.87	97.62
Chunking	97.80	97.41	96.52
Dependency UAS	94.52	94.13	93.59
Dependency LAS	92.61	92.16	91.47

Table 8: Effectiveness of using different layers for different tasks.

investigating such simple methods before developing complex methods for simple tasks. Currently, our JMT model does not explicitly use the learned dependency structures, and thus the explicit use of the output from the dependency layer should be an interesting direction of future work.

**Textual entailment** Table 6 shows the results of textual entailment, and our JMT model achieves the state-of-the-art result. The previous state-of-the-art result in Yin et al. (2016) relied on attention mechanisms and dataset-specific data pre-processing and features. Again, our simple max-pooling strategy achieves the state-of-the-art result boosted by the joint training. These results show the importance of jointly handling related tasks.

## 6.2 Analysis on the Model Architectures

We investigate the effectiveness of our model in detail. All of the results shown in this section are the development set results.

**Shortcut connections** Our JMT model feeds the word representations into all of the bi-LSTM layers, which is called the shortcut connection. Table 7 shows the results of “JMT<sub>all</sub>” with and without the shortcut connections. The results without the shortcut connections are shown in the column of “w/o SC”. These results clearly show that the importance of the shortcut connections, and in particular, the semantic tasks in the higher layers strongly rely on the shortcut connections. That is, simply stacking the LSTM layers is not sufficient to handle a variety of NLP tasks in a single model. In the supplementary material, it is qualitatively shown how the shortcut connections work in our model.

**Output label embeddings** Table 7 also shows the results without using the output labels of the POS, chunking, and relatedness layers, in the column of “w/o LE”. These results show that the explicit use of the output information from the classifiers of the lower layers is important in our JMT

	JMT <sub>all</sub>	w/o SR	w/o VC
POS	97.88	97.85	97.82
Chunking	97.59	97.13	97.45
Dependency UAS	94.51	94.46	94.38
Dependency LAS	92.60	92.57	92.48
Relatedness	0.236	0.239	0.241
Entailment	84.6	84.2	84.8

Table 9: Effectiveness of the Successive Regularization (SR) and the Vertical Connections (VC).

	JMT <sub>all</sub>	Random
POS	97.88	97.83
Chunking	97.59	97.71
Dependency UAS	94.51	94.66
Dependency LAS	92.60	92.80
Relatedness	0.236	0.298
Entailment	84.6	83.2

Table 10: Effects of the order of training.

model. The results in the column of “w/o SC&LE” are the ones without both of the shortcut connections and the label embeddings.

**Different layers for different tasks** Table 8 shows the results of our “JMT<sub>ABC</sub>” setting and that of not using the shortcut connections and the label embeddings (“w/o SC&LE”) as in Table 7. In addition, in the column of “All-3”, we show the results of using the highest (i.e., the third) layer for all of the three tasks without any shortcut connections and label embeddings, and thus the two settings “w/o SC&LE” and “All-3” require exactly the same number of the model parameters. The “All-3” setting is similar to the multi-task model of Collobert et al. (2011) in that task-specific output layers are used but most of the model parameters are shared. The results show that using the same layers for the three different tasks hampers the effectiveness of our JMT model, and the design of the model is much more important than the number of the model parameters.

**Successive regularization** In Table 9, the column of “w/o SR” shows the results of omitting the successive regularization terms described in Section 3. We can see that the accuracy of chunking is improved by the successive regularization, while other results are not affected so much. The chunking dataset used here is relatively small compared with other low-level tasks, POS tagging and dependency parsing. Thus, these results suggest that the successive regularization is effective when dataset sizes are imbalanced.

**Vertical connections** We investigated our JMT results without using the vertical connections in

	Single	Single+
POS	97.52	
Chunking	95.65	96.08
Dependency UAS	93.38	93.88
Dependency LAS	91.37	91.83
Relatedness	0.239	0.665
Entailment	83.8	66.4

Table 11: Effects of depth for the *single* tasks.

Single	W&C	Only W
POS	97.52	96.26
Chunking	95.65	94.92
Dependency UAS	93.38	92.90
Dependency LAS	91.37	90.44

Table 12: Effects of the character embeddings.

the five-layer bi-LSTMs. More concretely, when constructing the input vectors  $g_t$ , we do not use the bi-LSTM hidden states of the previous layers. Table 9 also shows the JMT<sub>all</sub> results with and without the vertical connections. As shown in the column of “w/o VC”, we observed the competitive results. Therefore, in the target tasks used in our model, sharing the word representations and the output label embeddings is more effective than just stacking the bi-LSTM layers.

**Order of training** Our JMT model iterates the training process in the order described in Section 3. Our hypothesis is that it is important to start from the lower-level tasks and gradually move to the higher-level tasks. Table 10 shows the results of training our model by randomly shuffling the order of the tasks for each epoch in the column of “Random”. We see that the scores of the semantic tasks drop by the random strategy. In our preliminary experiments, we have found that constructing the mini-batch samples from different tasks also hampers the effectiveness of our model, which also supports our hypothesis.

**Depth** The single task settings shown in Table 1 are obtained by using single layer bi-LSTMs, but in our JMT model, the higher-level tasks use successively deeper layers. To investigate the gap between the different number of the layers for each task, we also show the results of using multi-layer bi-LSTMs for the single task settings, in the column of “Single+” in Table 11. More concretely, we use the same number of the layers with our JMT model; for example, three layers are used for dependency parsing, and five layers are used for textual entailment. As shown in these results, deeper layers do not always lead to better results, and the joint learning is more important than mak-



ing the models complex only for single tasks.

**Character  $n$ -gram embeddings** Finally, Table 12 shows the results for the three single tasks with and without the pre-trained character  $n$ -gram embeddings. The column of “W&C” corresponds to using both of the word and character  $n$ -gram embeddings, and that of “Only W” corresponds to using only the word embeddings. These results clearly show that jointly using the pre-trained word and character  $n$ -gram embeddings is helpful in improving the results. The pre-training of the character  $n$ -gram embeddings is also effective; for example, without the pre-training, the POS accuracy drops from 97.52% to 97.38% and the chunking accuracy drops from 95.65% to 95.14%.

### 6.3 Discussion

**Training strategies** In our JMT model, it is not obvious when to stop the training while trying to maximize the scores of all the five tasks. We focused on maximizing the accuracy of dependency parsing on the development data in our experiments. However, the sizes of the training data are different across the different tasks; for example, the semantic tasks include only 4,500 sentence pairs, and the dependency parsing dataset includes 39,832 sentences with word-level annotations. Thus, in general, dependency parsing requires more training epochs than the semantic tasks, but currently, our model trains all of the tasks for the same training epochs. The same strategy for decreasing the learning rate is also shared across all the different tasks, although our growing gradient clipping method described in Section 5.2 helps improve the results. Indeed, we observed that better scores of the semantic tasks can be achieved before the accuracy of dependency parsing reaches the best score. Developing a method for achieving the best scores for all of the tasks at the same time is important future work.

**More tasks** Our JMT model has the potential of handling more tasks than the five tasks used in our experiments; examples include entity detection and relation extraction as in Miwa and Bansal (2016) as well as language modeling (Godwin et al., 2016). It is also a promising direction to train each task for multiple domains by focusing on domain adaptation (Søgaard and Goldberg, 2016). In particular, incorporating language modeling tasks provides an opportunity to use large text data. Such large text data was used in our

experiments to pre-train the word and character  $n$ -gram embeddings. However, it would be preferable to efficiently use it for improving the entire model.

**Task-oriented learning of low-level tasks** Each task in our JMT model is supervised by its corresponding dataset. However, it would be possible to learn low-level tasks by optimizing high-level tasks, because the model parameters of the low-level tasks can be directly modified by learning the high-level tasks. One example has already been presented in Hashimoto and Tsuruoka (2017), where our JMT model is extended to learning task-oriented latent graph structures of sentences by training our dependency parsing component according to a neural machine translation objective.

## 7 Conclusion

We presented a joint many-task model to handle multiple NLP tasks with growing depth in a single end-to-end model. Our model is successively trained by considering linguistic hierarchies, directly feeding word representations into all layers, explicitly using low-level predictions, and applying successive regularization. In experiments on five NLP tasks, our single model achieves the state-of-the-art or competitive results on chunking, dependency parsing, semantic relatedness, and textual entailment.

### Acknowledgments

We thank the anonymous reviewers and the Salesforce Research team members for their fruitful comments and discussions.

## References

- 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.
- Daniel Andor, Chris Alberti, David Weiss, Aliaksei Severyn, Alessandro Presta, Kuzman Ganchev, Slav Petrov, and Michael Collins. 2016. Globally Normalized Transition-Based Neural Networks. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2442–2452.
- Giuseppe Attardi and Felice Dell’Orletta. 2008. Chunking and Dependency Parsing. In *Proceedings of LREC 2008 Workshop on Partial Parsing*.

- Bernd Bohnet. 2010. Top Accuracy and Fast Dependency Parsing is not a Contradiction. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 89–97.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, and Hui Jiang. 2016. Enhancing and Combining Sequential and Tree LSTM for Natural Language Inference. *arXiv*, cs.CL 1609.06038.
- Do Kook Choe and Eugene Charniak. 2016. Parsing as Language Modeling. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2331–2336.
- Ronan Collobert, Jason Weston, Leon Bottou, Michael Karlen nad Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural Language Processing (Almost) from Scratch. *Journal of Machine Learning Research*, 12:2493–2537.
- Timothy Dozat and Christopher D. Manning. 2017. Deep Biaffine Attention for Neural Dependency Parsing. In *Proceedings of the 5th International Conference on Learning Representations*.
- 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 and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 334–343.
- Jason Eisner. 1996. Efficient Normal-Form Parsing for Combinatory Categorical Grammar. In *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, pages 79–86.
- Akiko Eriguchi, Kazuma Hashimoto, and Yoshimasa Tsuruoka. 2016. Tree-to-Sequence Attentional Neural Machine Translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 823–833.
- Jonathan Godwin, Pontus Stenetorp, and Sebastian Riedel. 2016. Deep Semi-Supervised Learning with Linguistically Motivated Sequence Labeling Task Hierarchies. *arXiv*, cs.CL 1612.09113.
- Ian J. Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. 2013. Max-out Networks. In *Proceedings of The 30th International Conference on Machine Learning*, pages 1319–1327.
- Alex Graves and Jurgen Schmidhuber. 2005. Frame-wise Phoneme Classification with Bidirectional LSTM and Other Neural Network Architectures. *Neural Networks*, 18(5):602–610.
- Kazuma Hashimoto and Yoshimasa Tsuruoka. 2017. Neural Machine Translation with Source-Side Latent Graph Parsing. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. To appear.
- Sepp Hochreiter and Jurgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Eliyahu Kiperwasser and Yoav Goldberg. 2016. Easy-First Dependency Parsing with Hierarchical Tree LSTMs. *Transactions of the Association for Computational Linguistics*, 4:445–461.
- Taku Kudo and Yuji Matsumoto. 2001. Chunking with Support Vector Machines. In *Proceedings of the Second Meeting of the North American Chapter of the Association for Computational Linguistics*.
- Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. 2016. Ask Me Anything: Dynamic Memory Networks for Natural Language Processing. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 1378–1387.
- Adhiguna Kuncoro, Miguel Ballesteros, Lingpeng Kong, Chris Dyer, Graham Neubig, and Noah A. Smith. 2017. What Do Recurrent Neural Network Grammars Learn About Syntax? In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1249–1258.
- Alice Lai and Julia Hockenmaier. 2014. Illinois-LH: A Denotational and Distributional Approach to Semantics. In *Proceedings of the 8th International Workshop on Semantic Evaluation*, pages 329–334.
- Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. 2015. Recurrent Convolutional Neural Networks for Text Classification. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 2267–2273.
- Zhizhong Li and Derek Hoiem. 2016. Learning without Forgetting. *CoRR*, abs/1606.09282.
- Wang Ling, Chris Dyer, Alan W Black, Isabel Trancoso, Ramon Fernandez, 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.
- Minh-Thang Luong, Ilya Sutskever, Quoc V. Le, Oriol Vinyals, and Lukasz Kaiser. 2016. Multi-task Sequence to Sequence Learning. In *Proceedings of the 4th International Conference on Learning Representations*.

- Xuezhe Ma and Eduard Hovy. 2016. End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1064–1074.
- Marco Marelli, Luisa Bentivogli, Marco Baroni, Raffaella Bernardi, Stefano Menini, and Roberto Zamparelli. 2014. SemEval-2014 Task 1: Evaluation of Compositional Distributional Semantic Models on Full Sentences through Semantic Relatedness and Textual Entailment. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 1–8.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26*, pages 3111–3119.
- Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. 2016. Cross-stitch Networks for Multi-task Learning. *CoRR*, abs/1604.03539.
- Makoto Miwa and Mohit Bansal. 2016. End-to-End Relation Extraction using LSTMs on Sequences and Tree Structures. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1105–1116.
- Yasumasa Miyamoto and Kyunghyun Cho. 2016. Gated Word-Character Recurrent Language Model. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1992–1997.
- Masataka Ono, Makoto Miwa, and Yutaka Sasaki. 2015. Word Embedding-based Antonym Detection using Thesauri and Distributional Information. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 984–989.
- Vu Pham, Theodore Bluche, Christopher Kermorvant, and Jerome Louradour. 2014. Dropout improves Recurrent Neural Networks for Handwriting Recognition. *CoRR*, abs/1312.4569.
- Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. 2016. Progressive Neural Networks. *CoRR*, abs/1606.04671.
- Anders Søgaard. 2011. Semi-supervised condensed nearest neighbor for part-of-speech tagging. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 48–52.
- Anders Søgaard and Yoav Goldberg. 2016. Deep multi-task learning with low level tasks supervised at lower layers. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 231–235.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15:1929–1958.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems 27*, pages 3104–3112.
- Jun Suzuki and Hideki Isozaki. 2008. Semi-Supervised Sequential Labeling and Segmentation Using Giga-Word Scale Unlabeled Data. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 665–673.
- Kai Sheng Tai, Richard Socher, and Christopher D. Manning. 2015. Improved Semantic Representations From Tree-Structured Long Short-Term Memory Networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1556–1566.
- Kristina Toutanova, Dan Klein, Christopher D Manning, and Yoram Singer. 2003. Feature-Rich Part-of-Speech Tagging with a Cyclic Dependency Network. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 173–180.
- Yoshimasa Tsuruoka, Yusuke Miyao, and Jun’ichi Kazama. 2011. Learning with Lookahead: Can History-Based Models Rival Globally Optimized Models? In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*, pages 238–246.
- 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 and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 323–333.
- John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2016. CHARAGRAM: Embedding Words and Sentences via Character n-grams. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1504–1515.
- Wenpeng Yin, Hinrich Schtze, Bing Xiang, and Bowen Zhou. 2016. ABCNN: Attention-Based Convolutional Neural Network for Modeling Sentence Pairs. *Transactions of the Association for Computational Linguistics*, 4:259–272.

Xingxing Zhang, Jianpeng Cheng, and Mirella Lapata. 2017. Dependency Parsing as Head Selection. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pages 665–676.

Yuan Zhang and David Weiss. 2016. Stack-propagation: Improved Representation Learning for Syntax. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1557–1566.

Yao Zhou, Cong Liu, and Yan Pan. 2016. Modelling Sentence Pairs with Tree-structured Attentive Encoder. In *Proceedings of the 26th International Conference on Computational Linguistics*, pages 2912–2922.

## Supplemental Material

### A Training Details

**Pre-training embeddings** We used the `word2vec` toolkit to pre-train the word embeddings. We created our training corpus by selecting lowercased English Wikipedia text and obtained 100-dimensional Skip-gram word embeddings trained with the context window size 1, the negative sampling method (15 negative samples), and the sub-sampling method ( $10^{-5}$  of the sub-sampling coefficient). We also pre-trained the character  $n$ -gram embeddings using the same parameter settings with the case-sensitive Wikipedia text. We trained the character  $n$ -gram embeddings for  $n = 1, 2, 3, 4$  in the pre-training step.

**Embedding initialization** We used the pre-trained word embeddings to initialize the word embeddings, and the word vocabulary was built based on the training data of the five tasks. All words in the training data were included in the word vocabulary, and we employed the *word-dropout* method (Kiperwasser and Goldberg, 2016) to train the word embedding for the unknown words. We also built the character  $n$ -gram vocabulary for  $n = 2, 3, 4$ , following Wieting et al. (2016), and the character  $n$ -gram embeddings were initialized with the pre-trained embeddings. All of the label embeddings were initialized with uniform random values in  $[-\sqrt{6}/(\dim + C), \sqrt{6}/(\dim + C)]$ , where  $\dim = 100$  is the dimensionality of the label embeddings and  $C$  is the number of labels.

**Weight initialization** The dimensionality of the hidden layers in the bi-LSTMs was set to 100. We

initialized all of the softmax parameters and bias vectors, except for the forget biases in the LSTMs, with zeros, and the weight matrix  $W_d$  and the root node vector  $r$  for dependency parsing were also initialized with zeros. All of the forget biases were initialized with ones. The other weight matrices were initialized with uniform random values in  $[-\sqrt{6}/(\text{row} + \text{col}), \sqrt{6}/(\text{row} + \text{col})]$ , where  $\text{row}$  and  $\text{col}$  are the number of rows and columns of the matrices, respectively.

**Optimization** At each epoch, we trained our model in the order of POS tagging, chunking, dependency parsing, semantic relatedness, and textual entailment. We used mini-batch stochastic gradient descent to train our model. The mini-batch size was set to 25 for POS tagging, chunking, and the SICK tasks, and 15 for dependency parsing. We used a gradient clipping strategy with growing clipping values for the different tasks; concretely, we employed the simple function:  $\min(3.0, \text{depth})$ , where  $\text{depth}$  is the number of bi-LSTM layers involved in each task, and 3.0 is the maximum value. The learning rate at the  $k$ -th epoch was set to  $\frac{\varepsilon}{1.0 + \rho(k-1)}$ , where  $\varepsilon$  is the initial learning rate, and  $\rho$  is the hyperparameter to decrease the learning rate. We set  $\varepsilon$  to 1.0 and  $\rho$  to 0.3. At each epoch, the same learning rate was shared across all of the tasks.

**Regularization** We set the regularization coefficient to  $10^{-6}$  for the LSTM weight matrices,  $10^{-5}$  for the weight matrices in the classifiers, and  $10^{-3}$  for the successive regularization term excluding the classifier parameters of the lower-level tasks, respectively. The successive regularization coefficient for the classifier parameters was set to  $10^{-2}$ . We also used *dropout* (Srivastava et al., 2014). The dropout rate was set to 0.2 for the vertical connections in the multi-layer bi-LSTMs (Pham et al., 2014), the word representations and the label embeddings of the entailment layer, and the classifier of the POS tagging, chunking, dependency parsing, and entailment. A different dropout rate of 0.4 was used for the word representations and the label embeddings of the POS, chunking, and dependency layers, and the classifier of the relatedness layer.



## B Details of Character $N$ -Gram Embeddings

Here we first describe the pre-training process of the character  $n$ -gram embeddings in detail and then show further analysis on the results in Table 12.

### B.1 Pre-Training with Skip-Gram Objective

We pre-train the character  $n$ -gram embeddings using the objective function of the Skip-gram model with negative sampling (Mikolov et al., 2013). We build the vocabulary of the character  $n$ -grams based on the training corpus, the case-sensitive English Wikipedia text. This is because such case-sensitive information is important in handling some types of words like named entities. Assuming that a word  $w$  has its corresponding  $K$  character  $n$ -grams  $\{cn_1, cn_2, \dots, cn_K\}$ , where any overlaps and unknown ones are removed. Then the word  $w$  is represented with an embedding  $v_c(w)$  computed as follows:

$$v_c(w) = \frac{1}{K} \sum_{i=1}^K v(cn_i), \quad (11)$$

where  $v(cn_i)$  is the parameterized embedding of the character  $n$ -gram  $cn_i$ , and the computation of  $v_c(w)$  is exactly the same as the one used in our JMT model explained in Section 2.1.

The remaining part of the pre-training process is the same as the original Skip-gram model. For each word-context pair  $(w, \bar{w})$  in the training corpus,  $N$  negative context words are sampled, and the objective function is defined as follows:

$$\sum_{(w, \bar{w})} \left( -\log \sigma(v_c(w) \cdot \tilde{v}(\bar{w})) - \sum_{i=1}^N \log \sigma(-v_c(w) \cdot \tilde{v}(\bar{w}_i)) \right), \quad (12)$$

where  $\sigma(\cdot)$  is the logistic sigmoid function,  $\tilde{v}(\bar{w})$  is the weight vector for the context word  $\bar{w}$ , and  $\bar{w}_i$  is a negative sample. It should be noted that the weight vectors for the context words are parameterized for the words without any character information.

### B.2 Effectiveness on Unknown Words

One expectation from the use of the character  $n$ -gram embeddings is to better handle unknown words. We verified this assumption in the single

task setting for POS tagging, based on the results reported in Table 12. Table 13 shows that the joint use of the word and character  $n$ -gram embeddings improves the score by about 19% in terms of the accuracy for unknown words.

We also show the results of the single task setting for dependency parsing in Table 14. Again, we can see that using the character-level information is effective, and in particular, the improvement of the LAS score is large. These results suggest that it is better to use not only the word embeddings, but also the character  $n$ -gram embeddings by default. Recently, the joint use of word and character information has proven to be effective in language modeling (Miyamoto and Cho, 2016), but just using the simple character  $n$ -gram embeddings is fast and also effective.

## C Analysis on Dependency Parsing

Our dependency parser is based on the idea of predicting a head (or parent) for each word, and thus the parsing results do not always lead to correct trees. To inspect this aspect, we checked the parsing results on the development set (1,700 sentences), using the “JMT<sub>ABC</sub>” setting.

In the dependency annotations used in this work, each sentence has only one root node, and we have found 11 sentences with multiple root nodes and 11 sentences with no root nodes in our parsing results. We show two examples below:

- (a) Underneath the headline “ Diversification , ” it **counsels** , “ Based on the events of the past week , all investors **need** to know their portfolios are balanced to help protect them against the market ’s volatility . ”
- (b) Mr. Eskandarian , who resigned his Della Femina post in September , becomes chairman and chief executive of Arnold .

In the example (a), the two boldfaced words “counsels” and “need” are predicted as child nodes of the root node, and the underlined word “counsels” is the correct one based on the gold annotations. This example sentence (a) consists of multiple internal sentences, and our parser misunderstood that both of the two verbs are the heads of the sentence.

In the example (b), none of the words is connected to the root node, and the correct child node of the root is the underlined word “chairman”.

Single (POS)	Overall Acc.	Acc. for unknown words
W&C	97.52	90.68 (3,502/3,862)
Only W	96.26	71.44 (2,759/3,862)

Table 13: POS tagging scores on the development set with and without the character  $n$ -gram embeddings, focusing on accuracy for unknown words. The overall accuracy scores are taken from Table 12. There are 3,862 unknown words in the sentences of the development set.

Single (Dependency)	Overall scores		Scores for unknown words	
	UAS	LAS	UAS	LAS
W&C	93.38	91.37	92.21 (900/976)	87.81 (857/976)
Only W	92.90	90.44	91.39 (892/976)	81.05 (791/976)

Table 14: Dependency parsing scores on the development set with and without the character  $n$ -gram embeddings, focusing on UAS and LAS for unknown words. The overall scores are taken from Table 12. There are 976 unknown words in the sentences of the development set.

Without the internal phrase “who resigned... in September”, the example sentence (b) is very simple, but we have found that such a simplified sentence is still not parsed correctly. In many cases, verbs are linked to the root nodes, but sometimes other types of words like nouns can be the candidates. In our model, the single parameterized vector  $r$  is used to represent the root node for each sentence. Therefore, the results of the examples (a) and (b) suggest that it would be needed to capture various types of root nodes, and using sentence-dependent root representations would lead to better results in future work.

## D Analysis on Semantic Tasks

We inspected the development set results on the semantic tasks using the “JMT<sub>all</sub>” setting. In our model, the highest-level task is the textual entailment task. We show an example premise-hypothesis pair which is misclassified in our results:

Premise: “A surfer is riding a *big* wave across dark green water”, and

Hypothesis: “The surfer is riding a *small* wave”.

The predicted label is `entailment`, but the gold label is `contradiction`. This example is very easy by focusing on the difference between the two words “big” and “small”. However, our model fails to correctly classify this example because there are few opportunities to learn the difference. Our model relies on the pre-trained word

embeddings based on word co-occurrence statistics (Mikolov et al., 2013), and it is widely known that such co-occurrence-based embeddings can rarely discriminate between antonyms and synonyms (Ono et al., 2015). Moreover, the other four tasks in our JMT model do not explicitly provide the opportunities to learn such semantic aspects. Even in the training data of the textual entailment task, we can find only one example to learn the difference between the two words, which is not enough to obtain generalization capacities. Therefore, it is worth investigating the explicit use of external dictionaries or the use of pre-trained word embeddings learned with such dictionaries (Ono et al., 2015), to see whether our JMT model is further improved not only for the semantic tasks, but also for the low-level tasks.

## E How Do Shared Embeddings Change

In our JMT model, the word and character  $n$ -gram embedding matrices are shared across all of the five different tasks. To better qualitatively explain the importance of the shortcut connections shown in Table 7, we inspected how the shared embeddings change when fed into the different bi-LSTM layers. More concretely, we checked closest neighbors in terms of the cosine similarity for the word representations before and after fed into the forward LSTM layers. In particular, we used the corresponding part of  $W_u$  in Eq. (1) to perform linear transformation of the input embeddings, because  $u_t$  directly affects the hidden states of the LSTMs. Thus, this is a context-independent analysis.

Table 15 shows the examples of the word “standing”. The row of “Embedding” shows the cases of using the shared embeddings, and the others show the results of using the linear-transformed embeddings. In the column of “Only word”, the results of using only the word embeddings are shown. The closest neighbors in the case of “Embedding” capture the semantic similarity, but after fed into the POS layer, the semantic similarity is almost washed out. This is not surprising because it is sufficient to cluster the words of the same POS tags: here, NN, VBG, etc. In the chunking layer, the similarity in terms of verbs is captured, and this is because it is sufficient to identify the coarse chunking tags: here, VP. In the dependency layer, the closest neighbors are adverbs, gerunds of verbs, and nouns, and all of them can be child nodes of verbs in dependency trees. However, this information is not sufficient in further classifying the dependency labels. Then we can see that in the column of “Word and char”, jointly using the character  $n$ -gram embeddings adds the morphological information, and as shown in Table 12, the LAS score is substantially improved.

In the case of semantic tasks, the projected embeddings capture not only syntactic, but also semantic similarities. These results show that different tasks need different aspects of the word similarities, and our JMT model efficiently transforms the shared embeddings for the different tasks by the simple linear transformation. Therefore, without the shortcut connections, the information about the word representations are fed into the semantic tasks after transformed in the lower layers where the semantic similarities are not always important. Indeed, the results of the semantic tasks are very poor without the shortcut connections.

	Word and char	Only word
Embedding	leaning kneeling saluting clinging railing	stood stands sit pillar cross-legged
POS	warning waxing dunking proving tipping	ladder rc6280 bethle warning f-a-18
Chunking	applauding disdaining pickin readjusting reclaiming	fight favor pick rejoin answer
Dependency	guaranteeing resting grounding hanging hugging	patiently hugging anxiously resting disappointment
Relatedness	stood stands unchallenged notwithstanding judging	stood unchallenged stands beside exists
Entailment	nudging skirting straddling contesting footing	beside stands pillar swung ovation

Table 15: Closest neighbors of the word “standing” in the embedding space and the projected space in each forward LSTM.