

# Unsupervised Pre-training with Structured Knowledge for Improving Natural Language Inference

Xiaoyu Yang Xiaodan Zhu Zhan Shi Tianda Li

ECE Department, Queen’s University, Canada

xiaoyu.yang@queensu.ca

## Abstract

While recent research on natural language inference has considerably benefited from large annotated datasets (Williams et al., 2017; Bowman et al., 2015), the amount of inference-related knowledge (including commonsense) provided in the annotated data is still rather limited. There have been two lines of approaches that can be used to further address the limitation: (1) unsupervised pretraining can leverage knowledge in much larger unstructured text data; (2) structured (often human-curated) knowledge has started to be considered in neural-network-based models for NLI. An immediate question is whether these two approaches complement each other, or how to develop models that can bring together their advantages. In this paper, we propose models that leverage structured knowledge in different components of pre-trained models. Our results show that the proposed models perform better than previous BERT-based state-of-the-art models. Although our models are proposed for NLI, they can be easily extended to other sentence or sentence-pair classification problems.

## 1 Introduction

Natural language inference (NLI), also known as recognizing textual entailment (RTE) (Bowman et al., 2015; Dagan et al., 2013; MacCartney and Manning, 2008), is a challenging problem in natural language understanding. It also acts as a test bed for representation learning for natural language (Bowman et al., 2015; Williams et al., 2017; Wang et al., 2018a). Specifically, NLI asks a system to identify the relationship between a premise  $p$  and a hypothesis  $h$  such as entailment and contradiction.

Recent advance on NLI has benefited from the availability of large annotated data (Williams et al., 2017; Bowman et al., 2015). It is, however, unlikely that the annotated data will contain all needed inference knowledge.

The most recent year has seen two lines of approaches that can be used to help address the limitation. First, the state-of-the-art unsupervised pre-training has shown to be very effective in leveraging knowledge in large unstructured data (Devlin et al., 2018; Radford et al., 2018; Peters et al., 2018). In parallel, some research has started to incorporate structured knowledge (Chen et al., 2018) into neural-network-based NLI models. Whether these two approaches complement each other in leveraging knowledge in large unstructured text and structured knowledge. How to develop models that can bring together the advantages of these two approaches?

In this paper, we explore a variety of approaches to leverage structured external knowledge in Transformer-based pre-training architectures (Vaswani et al., 2017), which are widely used in the state-of-the-art pre-trained models such as generative pre-trained Transformer (GPT) model (Radford et al., 2018) and BERT (Devlin et al., 2018). Specifically we (1) incorporate structured knowledge in Transformer self-attention, (2) add a knowledge-specific layer over the existing Transformer block, and (3) incorporate structured knowledge in global inference. We believe that complicated NLP tasks such as NLI can benefit from knowledge available outside the training data from these two typical sources: the knowledge learned implicitly from unstructured text through unsupervised pre-training and that from a structured (often human-created) knowledge base.

We evaluate the proposed models on the widely used NLI benchmarks, including the Stanford Natural Language Inference (SNLI) corpus (Bowman et al., 2015), Multi-Genre Natural Language Inference (MultiNLI) corpus (Williams et al., 2017), and a newly introduced diagnostic Glockner dataset (Glockner et al., 2018). Our results show that the proposed models perform better on these datasets than previous BERT-based state-of-the-art

models. While our models are proposed for NLI, they can be easily adapted to other sentence or sentence-pair classification problems.

## 2 Related Work

### 2.1 Pre-training in NLI and Other Tasks

Unsupervised pre-trained models (Peters et al., 2018; Radford et al., 2018; Devlin et al., 2018; Radford et al., 2019; Song et al., 2019) have proven to be effective in achieving great improvements in many NLP tasks. Feature-based and finetune-based are two strategies where the pre-trained models apply to downstream tasks. Feature-based models such as ELMo (Peters et al., 2018) provide pre-trained representations as additional features to original task-specific models. Finetune-based models such as Generative pre-trained Transformer (GPT) (Radford et al., 2018) and BERT (Devlin et al., 2018), however, use pre-trained architectures and then fine-tune the same architectures on downstream tasks. In this paper, we explore combining structured knowledge into pre-trained models in the stage of fine-tuning.

### 2.2 Human-Authorized External Knowledge Integration

As discussed above, while neural networks have been shown to be very effective in modeling NLI with large amounts of training data, their end-to-end training is based on a strong assumption that all necessary knowledge for inference is learnable from the provided training data. This kind of assumption, however, can be limited in some circumstances where the knowledge available in a human-authorized external knowledge base is too sparse to learn just from the training data or corpora the model is pre-trained with.

Existing work has made an effort to incorporate human-authorized knowledge (e.g., WordNet, ConceptNet) in neural networks. For entity and relation representation of knowledge graph, there are both structure-based and semantically-enriched approaches to get their embeddings. For example, TransE (Bordes et al., 2011) provides structure-based embedding by modeling relations as translations, while for Neural Tensor Network (Socher et al., 2013), it represents entities as an average of their word vectors and then it is used for tensor-based transformation. After getting vector representation for entities or relations, they can be used to optimize sentence alignment like (Chen

et al., 2018). Some efforts have also been made to generate knowledge enhanced representations for premise-hypothesis pairs and then used to benefit the final classification like (Wang et al., 2018b).

### 2.3 Evaluation of Models for NLI

Previous research has made contributions to large annotated datasets for NLI evaluation. For example, SNLI (Bowman et al., 2015) and MultiNLI (Williams et al., 2017) are two widely used NLI datasets. Previous research (Gururangan et al., 2018; Naik et al., 2018) has also paid attention to whether existing NLI systems have learned NLI-related semantics or just explored the regularities existing in the data that are not relevant to NLI.

To investigate this, different methods have been proposed. (Wang et al., 2019) introduces a *swapping* evaluation method, which means changing the distribution of words by swapping a premise with its corresponding hypothesis to test the robustness of models. Also, new test datasets are proposed, e.g., Glockner test set (Glockner et al., 2018). In the Glockner test dataset, premises are taken from the SNLI training set, and hypotheses are generated by replacing a single word in its corresponding premise sentence. In addition, swapping Glockner test in (Li et al., 2019) shows that external knowledge from unsupervised pre-training and human-authorized external knowledge can benefit models by improving the capacity of capturing NLI semantics. Also, human-authorized external knowledge has the potential to enhance pre-trained models.

## 3 The Models

We propose three approaches to leverage structured external knowledge in Transformer-based pre-training frameworks. Figure 1 shows an overview of our models.

- **Structured Knowledge for Attention Weights Adjustment** As shown in the blue block (M1) of Fig 1, this approach incorporates structured knowledge in multi-head self-attention layer of Transformer.
- **Structured Knowledge for a Separate Knowledge-Specific Layer** As shown in the yellow block (M2) of Fig 1, this approach adds an additional knowledge-specific multi-head self-attention layer in Transformer.

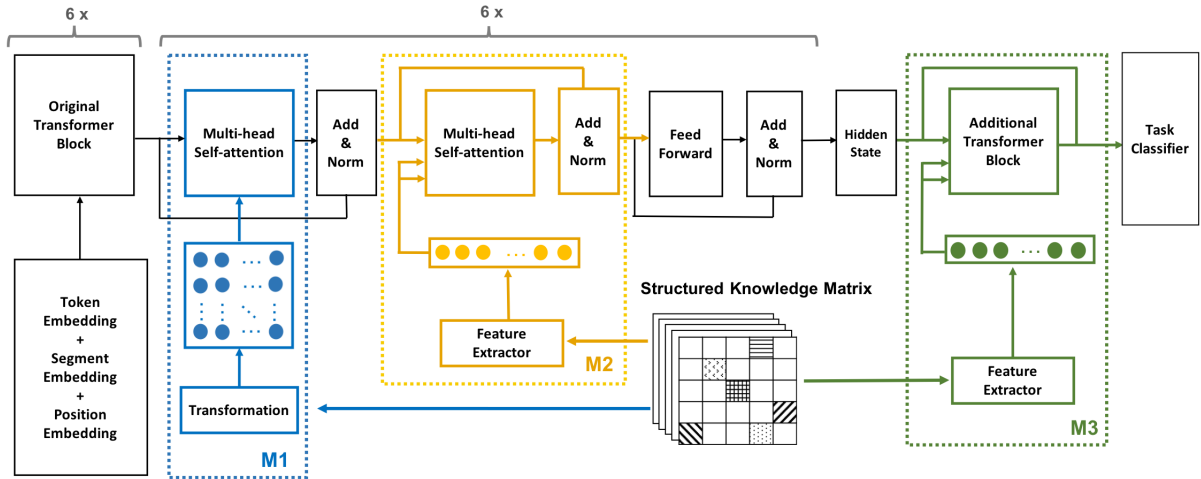


Figure 1: An overview of the three different approaches incorporating structured knowledge on  $BERT_{BASE}$  ( $12 \times$  Transformer). The blue, yellow and green frames denote approaches M1, M2 and M3 respectively.

- **Structured Knowledge for Global Inference** As shown in in the green block (M3), this approach directly incorporates structured knowledge into global inference.

We implement all our models based on  $BERT_{BASE}$ <sup>1</sup>, and our approaches could be easily extended to other Transformer-based architectures such as GPT.

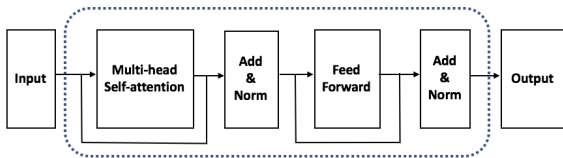


Figure 2: Transformer architecture.

To make the content complete and help introduce our models, we briefly describe Transformer here, and the details can be found in (Vaswani et al., 2017). Transformer is widely used as a building block of many state-of-the-art NLP models. Figure 2 shows the architecture of Transformer, which contains a multi-head self-attention layer and a feed-forward layer. Residual connections and normalization are applied between layers.  $BERT_{BASE}$  is built with twelve identical Transformer blocks. The network is randomly initialized and then pre-trained on BooksCorpus (Zhu et al., 2015) and a Wikipedia Corpus with two objectives, *i.e.*, masked LM and next sentence prediction (Devlin et al., 2018).

<sup>1</sup><https://github.com/google-research/bert>

### 3.1 Structured Knowledge

In general, structured knowledge in relational databases or knowledge graphs can be regarded as a list of triples  $\langle n_1, n_2, r \rangle$ , where  $n_1$  and  $n_2$  are nodes and  $r$  is their relation. In this work, we use triples from commonsense knowledge bases, WordNet (Miller, 1995) and ConceptNet 5.5 (Speer et al., 2017), as the sources of structured knowledge. Although there exist other sources of structured knowledge, *e.g.*, Freebase (WikiData) (Bollacker et al., 2007), such knowledge sources mainly convey factual evidence (*e.g.*, the relationship between Bill Gates and Microsoft), which is less relevant for general natural language inference. A common example of general NLI is a sentence pair like “a girl in a yellow dress with the sun shining on her face” and “a girl in a pink dress with the sun shining on her face”, and determining the relationship between these two sentences relies more on the knowledge of the color “yellow” and “pink” are two different colors.

Structured knowledge could benefit natural language inference in two aspects. First, external knowledge may in general help align inference-related concepts between a premise and the corresponding hypothesis. For example, the features like *synonymy*, *antonymy*, *hypernymy*, *hyponymy* and *co-hyponyms* can help determine local inference (Chen et al., 2017), by helping aligning the mentions of such words between a premise and a hypothesis through attention. Second, such semantic relations could be directly used as inference-related features, thus providing auxiliary infor-

mation for inference. For example, knowledge of *hypernymy* and *hyponymy* may help capture entailment, and knowledge of *antonymy* and *co-hyponyms* may help model contradiction. Considering a sentence pair like “a man is drinking beer” and “a man is drinking whisky”, with the knowledge that “beer” and “whisky” are different kinds of alcohol (*co-hyponyms*), we can infer that these sentences contradict each other.

In conclusion, given a pair of premise and hypothesis for inference, each word pair  $\langle w_i, w_j \rangle$  between the two sentences will be assigned with a multi-dimensional semantic relation  $e_{ij}$ ,  $e_{ij} \in \mathbb{R}^k$ , where  $k$  represents the dimension of extracted knowledge. Particularly, in the BERT architecture, the premise and hypothesis are concatenated into a sequence of length  $n$  (zero-padding shorter sequence or truncate longer sequence to  $n$ ). Therefore, we can form a knowledge matrix  $E = \{e_{ij}\}$  with the shape of  $n \times n \times k$  for the pair of sentences.

### 3.2 Structured Knowledge for Attention Weights Adjustment

As the structured knowledge can be viewed as a prior to align tokens between premise and hypothesis, we incorporate external knowledge in the multi-head self-attention module of Transformer in the BERT<sub>BASE</sub> architecture. The multi-head self-attention module in each Transformer block contains multiple parallel heads and attention weight matrix  $A \in \mathbb{R}^{n \times n}$  for each head is calculated by the pair-wise dot product of the transformed input representations, where  $n$  represents the length of the input sequence.

Specifically, we represent the input sequence as  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t, \dots, \mathbf{x}_n]$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  is a continuous representation for a token at position  $t$ , and  $x_t$  is an element-wise summation of token embedding, segment embedding, and position embedding. In each head, we first calculate query ( $Q$ ), key ( $K$ ) and Value ( $V$ ) as below, which are different transformations of input  $X$ .

$$Q = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t, \dots, \mathbf{x}_n]W_q \quad (1)$$

$$K = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t, \dots, \mathbf{x}_n]W_k \quad (2)$$

$$V = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t, \dots, \mathbf{x}_n]W_v \quad (3)$$

The output  $H$  as well as the attention weight

matrix  $A$  are computed as follows:

$$A = \{a_{ij}\} = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \quad (4)$$

$$H = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_t, \dots, \mathbf{h}_n]^T = AV \quad (5)$$

where as in Transformer (Vaswani et al., 2017),  $\sqrt{d_k}$  acts as a scaling factor when calculating the attention weight matrix  $A$ .

We combine structured knowledge  $E \in \mathbb{R}^{n \times n \times k}$  with self-attention weight  $A \in \mathbb{R}^{n \times n}$  for each head in the top six Transformer blocks while keeping the bottom six Transformer blocks unchanged. As shown in M1 block of Figure 1, the structured knowledge matrix is firstly fed into a transformation block to transform into a new matrix  $E' = \{e'_{ij}\}$  with the shape  $n \times n$  (We choose the average pooling along the last dimension of  $E$  for the transformation block). And then we merge  $E'$  into  $A$  to get a new attention weight matrix  $A'$  as follows:

$$A' = A + A \odot E' = \{a_{ij} + a_{ij}e'_{ij}\} \quad (6)$$

Here  $A'$  has the same shape as  $A$ , and can be directly used to replace  $A$  in Equation 5 to get the structured knowledge augmented output  $H$ . In this way, we inject structured knowledge in each multi-head self-attention layers so as to improve the alignments between tokens.

### 3.3 Structured Knowledge for a Separate Knowledge-Specific Layer

Intuitively, the structured knowledge retrieved from human-curated knowledge bases could contain helpful information to determine the relationship between premise and hypothesis. Therefore, we explore to encode the structured knowledge as knowledge features and then combine these features into the top six Transformer blocks by an additional multi-head self-attention layer shown in M2 of Figure 1. Also, residual connection and layer normalization are included in this additional layer.

The extractor used here for knowledge feature extracting is a Convolutional Neural Network (CNN) composed of a stacked convolutional layer with  $p$  different filters (Equation 7) as well as consecutive pooling layers which take its input as the concatenation of the output from convolutional layers (Equation 8). The output  $H$  of the previous self-attention layer in Transformer block is regarded as query ( $Q$ ) of this additional multi-head self-attention layer, while the extracted knowledge



features  $C$  are used as its corresponding key ( $K$ ) and value ( $V$ ). Specifically, the output of this newly introduced self-attention module is calculated as below:

$$c_i = \text{Conv}(E, \text{filter}_i), i \in \{1, \dots, p\} \quad (7)$$

$$C = \text{pooling}[c_1; c_2; \dots; c_p] \quad (8)$$

$$A' = \{a'_{ij}\} = \text{softmax}\left(\frac{HC^T}{\sqrt{d_k}}\right) \quad (9)$$

$$P = [p_1, p_2, \dots, p_n]^T = A'C \quad (10)$$

Where  $A'$  and  $P$  denote the attention weight and output in the additional multi-head self-attention layer respectively. In this way, we incorporate inference-related knowledge to benefit the reasoning.

### 3.4 Structured Knowledge for Global Inference

The structured knowledge matrix is expected to provide explicit semantic information to indicate the overall label for a pair of premise and hypothesis. We explore incorporating structured knowledge directly in the global inference layer by adding an additional Transformer block on top of the original BERT<sub>base</sub> model. This additional Transformer block is a simplified version of a standard Transformer block in BERT with one head instead of multiple ones, and its parameters are randomly initialized and optimized during training.

The structured knowledge matrix  $E$  is first fed into a feature extractor, which generates a structured knowledge matrix  $M$  according to Equation 11 to 12. Similar to the feature extractor used in Section 3.3, the feature extractor we use in this approach is also a Convolution Neural Network.

$$m_i = \text{Conv}(E, \text{filter}_i), i \in \{1, \dots, p'\} \quad (11)$$

$$M = \text{pooling}[m_1; m_2; \dots; m_{p'}] \quad (12)$$

$$w = \text{softmax}\left(\frac{Mh_0^T}{\sqrt{d_k}}\right) \quad (13)$$

$$h_{final} = Mw \quad (14)$$

Where  $h_0$  is the output hidden vector of the first token in output sequence as the compact sentence representation and  $h_{final}$  is then fed into the classification layer. As shown in M3 of Figure 1, for this approach, we combine the structured knowledge matrix  $M$  with the compact sentence representation vector  $h_0$  by designing an additional Transformer

block to produce the input  $h_{final}$  for the task classifier. In this way, we manage to leverage structured knowledge in global inference.

## 4 Experiment Setup

### 4.1 Data

We evaluate the proposed models on the widely used NLI datasets, *i.e.*, SNLI and MultiNLI (MultiNLI-match and MultiNLI-mismatch). We also test our models with the Glockner test-set (Glockner et al., 2018). These datasets share the same target of a 3-way prediction: determining the relation in a premise-hypothesis pair to be either *entailment*, *neutral*, or *contradiction*.

For SNLI and MultiNLI, our models are trained and tested based on official splits of training, development, and test set. For the Glockner test set, our models are trained on SNLI.

### 4.2 Representation of Structured External Knowledge

The most prominent structured external knowledge for downstream NLI task is lexicon semantics (phrase-level inference knowledge is much more sparse and hard to acquire). We use WordNet and ConceptNet as structured knowledge sources and extract specific semantic relations between word pairs. WordNet is a human-authorized knowledge base in which lexical semantics such as those about nouns, verbs, adjectives, and adverbs are organized into sets of synonyms with semantic relations linking them. ConceptNet5.5 includes massive world and lexical knowledge in different languages from different sources and organizes such knowledge in a graph with differently weighted edges.

In WordNet, we use five types of semantic relation, including *hypernymy*, *hyponymy*, *co-hyponyms*, *antonymy*, and *synonymy*. We extract and represent this knowledge with the method proposed in (Chen et al., 2017). The dimension  $k$  of semantic relation  $e_{ij}$  is 5 in our case, if a specific word pair falls into either synonymy or antonymy relation, then the corresponding dimension will be 1, otherwise it will be 0. If two words do not belong to the same synset but share the same hypernym, the value of co-hyponyms dimension will be set as 1, otherwise it is 0. When calculating hypernymy features, it takes the value  $(1 - n/8)$  if one word is a hypernym of the other word within 8 steps in the WordNet hierarchy. Given calculated hypernymy relation between word  $A$  and word  $B$ ,

WordNet	ConceptNet
Hypernymy	HasA
Hyponymy	InstanceOf, Entails, IsA, MannerOf, MadeOf, PartOf, DerivedFrom
Co-hyponyms	DistinctFrom
Antonymy	Antonym
Synonymy	FormOf, SimilarTo, Synonym

Table 1: Corresponding relation between WordNet and ConceptNet.

	WordNet	ConceptNet
Hypernymy	753,086	5,532
Hyponymy	753,086	434,381
Co-hyponyms	3,674,700	3,396
Antonymy	6,617	18,625
Synonymy	237,937	602,399

Table 2: Statistics for semantic relation pairs extracted from WordNet and ConceptNet.

like  $[A, B] = 0.125$ , it is easy to infer that the hyponymy feature is  $[B, A] = 0.125$ .

We also use ConceptNet 5.5 as a structured knowledge source. The knowledge from ConceptNet will only be added when there is no counterpart in WordNet. As mentioned before, the structured knowledge representation will be a 5-dimensional vector. In ConceptNet 5.5, the concepts are aligned into 36 relations. In order to import ConceptNet, we managed to condense ConceptNet relation features compatible with those from WordNet. Specifically, we filter out relations similar to those selected from WordNet and establish the corresponding relationship between relations from WordNet and ConceptNet as shown in Table 1.

Table 2 provides the statistics about how many semantic relations for word pairs can be extracted from WordNet and ConceptNet. For ConceptNet, only pairs with both words being English are counted.

### 4.3 Model Details

Our model is based on the official pre-trained  $BERT_{BASE}$ . During fine-tuning, we load all the pre-trained parameters from the officially released pre-trained  $BERT_{BASE}$  model, and then we initialize other parameters, including those for the task-specific classifiers and those we newly introduced in our model.  $BERT_{BASE}$  consists of 12

Transformer building blocks, in which the number of heads for the attention module is 12, and the size of the hidden state is 768. The maximum length of the input sequence is 128.

The feature extractor we used in Section 3.3 is CNN, with a 4-layer stacked convolutional layer and following two pooling layers. The filters for the stacked convolutional layer have kernel sizes of  $3 \times 3$ ,  $5 \times 5$ ,  $7 \times 7$ , and  $9 \times 9$ . The number of channels for each layer is 16 with a stride of 1. After feeding the same structured knowledge matrix into different convolutional layers, we concatenated the four outputs by the last dimension, and then it is followed by 2 max pooling layers, which have pooling sizes of  $2 \times 2$  and  $5 \times 5$ , and are applied at stride 2 and 3, respectively. There is a single-layer feed forward network at the end to transform the structured knowledge vector into an appropriate size for further processing.

The feature extractor we used in Section 3.4 is a CNN with a similar architecture as the CNN used in Section 3.3. One of the different points is that it includes a 3-layer stacked convolutional layer instead of a 4-layer stacked one, and the filters are with the shape of  $3 \times 3$ ,  $5 \times 5$ , and  $7 \times 7$ . The other difference is mainly about the parameters of the single-layer feed forward network since we need to transform the knowledge vector into different space with a different dimension.

## 5 Experiment Results

### 5.1 Overall Performance

Table 3 presents the results of our model on SNLI, MultiNLI-matched, MultiNLI-mismatched, as well as the Glockner test set. Our model here represents the overall model that combines M1, M2, and M3 discussed before. We first test our model with WordNet while without structured knowledge from ConceptNet. It shows that even using only WordNet can bring improvements compared with the baseline model. The best performance was achieved by leveraging knowledge from both WordNet and ConceptNet. Specifically, our model with the complete set of structured knowledge outperforms the baselines consistently across all these datasets compared with BERT. It shows that leveraging structured knowledge on pre-trained models can benefit NLI tasks. We have not experimented with  $BERT_{LARGE}$  because it is computationally expensive, but we believe the proposed model is orthogonal and useful, particularly when structured

Table 3: Results on NLI datasets, comparing our model with previous models. Our model is the combination of the three approaches proposed in this paper. The baseline results are from their published papers except those indicated with  $\star$  are based on our experiments with the officially released pre-trained models.

Model	SNLI	MultiNLI-m	MultiNLI-mm	SNLI-Glockner
ESIM (Chen et al., 2016)	88.0	76.8	75.8	65.6
KIM (Chen et al., 2018)	88.6	77.2	76.4	83.5
ESIM + ELMo (Peters et al., 2018)	89.3	-	-	-
GPT (Radford et al., 2018)	89.9	82.1	81.4	-
BERT <sub>BASE</sub> (Devlin et al., 2018)	90.6 $\star$	84.6	83.4	94.7 $\star$
Our Model (w/o ConceptNet)	90.9	84.7	84.1	<b>95.3</b>
Our Model (w ConceptNet)	<b>91.0</b>	<b>84.9</b>	<b>84.3</b>	<b>95.3</b>

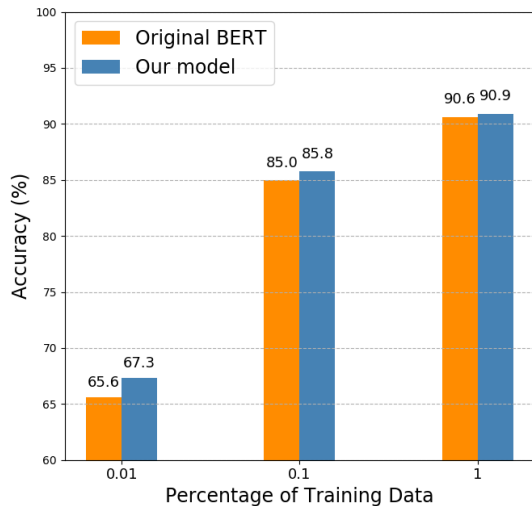


Figure 3: Accuracies of original BERT and our structured knowledge (with WordNet) enhanced model with different proportion of training data on SNLI.

knowledge is complementary to knowledge learned in pretraining, e.g., when domain-specific tasks are concerned with abundant structured domain knowledge available.

## 5.2 Ablation Analysis

To further prove that structured external knowledge can benefit BERT on NLI, we perform ablation tests to get empirical results for analysis.

The first test aims to observe the performance of M1, M2, and M3, respectively. The performances of three different approaches are shown in Table 4. Comparing with the results from BERT<sub>BASE</sub>, we can see that every single method can bring improvements on these datasets to different degrees.

In Figure 3, we present the test accuracy from

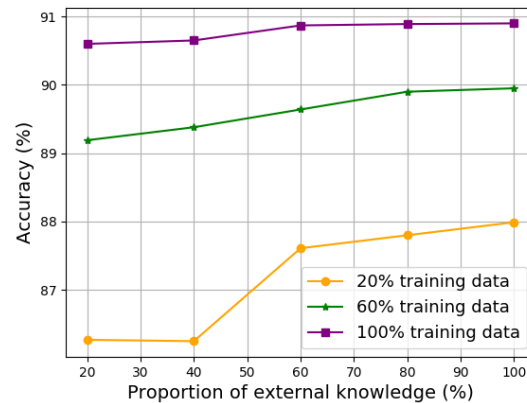


Figure 4: Accuracies of the overall model leveraging different proportion of structured knowledge (WordNet), and trained with different proportion of SNLI training data.

models trained with/without structured external knowledge. A comparison of both models under each data split shows that the model with structured knowledge performs consistently better than the original BERT baseline. The gap between the two models becomes larger when less training data is available. When we have only 0.01 percent of the data available, structured knowledge has more obvious effects on model performance. Our model shows extra benefit from the structured relation pairs, complementary to information from training set and pre-training.

Figure 4 shows increasing structured knowledge has a positive effect on test accuracy. The slope of each single line in Figure 4 shows that when we increase the amount of structured knowledge from 20% to 100%, the test accuracy increases accordingly, and the accuracy increases more obviously when there is less training data available. This in-

Table 4: Accuracies of M1, M2, and M3 on SNLI, MultiNLI-m, MultiNLI-mm, and Glockner. M1, M2 and M3 refer to the submodel in Section 3.2, 3.3, and 3.4, respectively. The resource of structured knowledge is WordNet, and the models tested on Glockner are trained with SNLI.

Model	SNLI	MultiNLI-m	MultiNLI-mm	SNLI-Glockner
BERT <sub>BASE</sub> (Devlin et al., 2018)	90.6	84.6	83.4	94.7
M1 (w/o ConceptNet)	90.7	<b>85.0</b>	84.1	95.0
M2 (w/o ConceptNet)	90.7	84.6	<b>84.2</b>	95.0
M3 (w/o ConceptNet)	<b>90.9</b>	<b>85.0</b>	83.9	<b>95.6</b>

creasing trend provides evidence that the proposed model is capable of using lexical-level structured knowledge pairs for natural language inference.

Premise-Hypothesis pairs	P / T
<i>p</i> : A man playing an electric <b>guitar</b> on stage. <i>h</i> : A man playing <b>banjo</b> on the floor.	E / C
<i>p</i> : A man juggles bowling pins for the <b>first</b> time. <i>h</i> : A man juggles bowling pins for the <b>1st</b> time.	C / E
<i>p</i> : A woman is sipping some <b>wine</b> . <i>h</i> : A woman is sipping some <b>vodka</b> .	N / C

Table 5: Examples from Glockner and SNLI on which our model predicts correctly while not the original BERT. P / T represents prediction and true label. E, N, and C refer to entailment, neutral and contradiction, respectively.

### 5.3 Case Study

To further find out how structured knowledge enhances the pre-trained BERT, we analyze some specific cases from SNLI and Glockner, on which only our model can yield correct predictions but not original BERT, as shown in Table 5. Note that all the highlighted keyword pairs can be retrieved in the structured knowledge base we use. Take the first pair of sentences in Table 5 as an example, the highlighted pair of words “guitar” and “banjo” in the first example can be retrieved in WordNet with the relation of *co-hyponyms*, which indicates they have the same hypernym “instrument” but refer to different entities. Our model can leverage this kind of knowledge to make a reasonable inference that these two sentences contradict to each other as this kind of lexical semantic knowledge is too sparse to learn from unstructured training data.

## 6 Conclusions and Discussion

Research on natural language inference has considerably benefited from large annotated datasets. Most recently, there have been two lines of approaches to employing knowledge available outside training data: knowledge from structured knowledge bases and that learned from unsupervised pre-training. In this paper, we explore whether these two approaches complement each other, and how to develop models that can bring together their advantages. We propose models that leverage structured knowledge in different components of pre-trained models. The results show that the proposed models perform better than the BERT baseline. While our models are proposed for NLI, they can be easily extended to other sentence or sentence-pair classification problems to leverage these two sources of external knowledge.

## References

- Kurt Bollacker, Robert Cook, and Patrick Tufts. 2007. Freebase: A shared database of structured general human knowledge. In *AAAI*, volume 7, pages 1962–1963.
- Antoine Bordes, Jason Weston, Ronan Collobert, and Yoshua Bengio. 2011. Learning structured embeddings of knowledge bases. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. *A large annotated corpus for learning natural language inference*. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Diana Inkpen, and Si Wei. 2017. Neural natural language inference models enhanced with external knowledge. *arXiv preprint arXiv:1711.04289*.
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Diana Inkpen, and Si Wei. 2018. Neural natural language



- inference models enhanced with external knowledge. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 2406–2417.
- Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2016. Enhanced lstm for natural language inference. *arXiv preprint arXiv:1609.06038*.
- Ido Dagan, Dan Roth, Mark Sammons, and Fabio Massimo Zanzotto. 2013. Recognizing textual entailment: Models and applications. *Synthesis Lectures on Human Language Technologies*, 6(4):1–220.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. Breaking nli systems with sentences that require simple lexical inferences. *arXiv preprint arXiv:1805.02266*.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. [Annotation artifacts in natural language inference data](#).
- Tianda Li, Xiaodan Zhu, Quan Liu, Qian Chen, Zhigang Chen, and Si Wei. 2019. Several experiments on investigating pretraining and knowledge-enhanced models for natural language inference. *arXiv preprint arXiv:1904.12104*.
- Bill MacCartney and Christopher D Manning. 2008. Modeling semantic containment and exclusion in natural language inference. In *Proceedings of the 22nd International Conference on Computational Linguistics-Volume 1*, pages 521–528. Association for Computational Linguistics.
- George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. *arXiv preprint arXiv:1806.00692*.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. URL [https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language\\_understanding\\_paper.pdf](https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language_understanding_paper.pdf).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Richard Socher, Danqi Chen, Christopher D Manning, and Andrew Ng. 2013. Reasoning with neural tensor networks for knowledge base completion. In *Advances in neural information processing systems*, pages 926–934.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tiejun Liu. 2019. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*.
- Robert Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Alex Wang, Amapreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018a. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Haohan Wang, Da Sun, and Eric P Xing. 2019. What if we simply swap the two text fragments? a straightforward yet effective way to test the robustness of methods to confounding signals in nature language inference tasks. 33(01):7136–7143.
- Xiaoyan Wang, Pavan Kapanipathi, Ryan Musa, Mo Yu, Kartik Talamadupula, Ibrahim Abdelaziz, Maria Chang, Achille Fokoue, Bassem Makni, Nicholas Mattei, et al. 2018b. Improving natural language inference using external knowledge in the science questions domain. *arXiv preprint arXiv:1809.05724*.
- Adina Williams, Nikita Nangia, and Samuel R Bowman. 2017. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27.