

GlossBERT: BERT for Word Sense Disambiguation with Gloss Knowledge

Luyao Huang, Chi Sun, Xipeng Qiu*, Xuanjing Huang

Shanghai Key Laboratory of Intelligent Information Processing, Fudan University

School of Computer Science, Fudan University

825 Zhangheng Road, Shanghai, China

{lyhuang18, sunc17, xpqiu, xjhuang}@fudan.edu.cn

Abstract

Word Sense Disambiguation (WSD) aims to find the exact sense of an ambiguous word in a particular context. Traditional supervised methods rarely take into consideration the lexical resources like WordNet, which are widely utilized in knowledge-based methods. Recent studies have shown the effectiveness of incorporating gloss (sense definition) into neural networks for WSD. However, compared with traditional *word expert* supervised methods, they have not achieved much improvement. In this paper, we focus on how to better leverage gloss knowledge in a supervised neural WSD system. We construct *context-gloss* pairs and propose three BERT-based models for WSD. We fine-tune the pre-trained BERT model on SemCor3.0 training corpus and the experimental results on several English all-words WSD benchmark datasets show that our approach outperforms the state-of-the-art systems¹.

1 Introduction

Word Sense Disambiguation (WSD) is a fundamental task and long-standing challenge in Natural Language Processing (NLP), which aims to find the exact sense of an ambiguous word in a particular context (Navigli, 2009). Previous WSD approaches can be grouped into two main categories: knowledge-based and supervised methods.

Knowledge-based WSD methods rely on lexical resources like WordNet (Miller, 1995) and usually exploit two kinds of lexical knowledge. The gloss, which defines a word sense meaning, is first utilized in Lesk algorithm (Lesk, 1986) and then widely taken into account in many other approaches (Banerjee and Pedersen, 2002; Basile et al., 2014). Besides, structural properties of semantic graphs are mainly used in graph-based algorithms (Agirre et al., 2014; Moro et al., 2014).

Traditional supervised WSD methods (Zhong and Ng, 2010; Shen et al., 2013; Iacobacci et al., 2016) focus on extracting manually designed features and then train a dedicated classifier (*word expert*) for every target lemma.

Although *word expert* supervised WSD methods perform better, they are less flexible than knowledge-based methods in the all-words WSD task (Raganato et al., 2017a). Recent neural-based methods are devoted to dealing with this problem. Kågebäck and Salomonsson (2016) present a supervised classifier based on Bi-LSTM, which shares parameters among all word types except the last layer. Raganato et al. (2017a) convert WSD task to a sequence labeling task, thus building a unified model for all polysemous words. However, neither of them can totally beat the best *word expert* supervised methods.

More recently, Luo et al. (2018b) propose to leverage the gloss information from WordNet and model the semantic relationship between the context and gloss in an improved memory network. Similarly, Luo et al. (2018a) introduce a (hierarchical) co-attention mechanism to generate co-dependent representations for the context and gloss. Their attempts prove that incorporating gloss knowledge into supervised WSD approach is helpful, but they still have not achieved much improvement, because they may not make full use of gloss knowledge.

In this paper, we focus on how to better leverage gloss information in a supervised neural WSD system. Recently, the pre-trained language models, such as ELMo (Peters et al., 2018) and BERT (Devlin et al., 2018), have shown their effectiveness to alleviate the effort of feature engineering. Especially, BERT has achieved excellent results in question answering (QA) and natural language inference (NLI). We construct *context-gloss* pairs

*Corresponding author.

¹<https://github.com/HSLCY/GlossBERT>

Sentence with four targets:

Your research stopped when a convenient assertion could be made.

Context-Gloss Pairs of the target word [research]

	Label	Sense Key
[CLS] Your research ... [SEP] systematic investigation to ... [SEP]	Yes	research%1:04:00::
[CLS] Your research ... [SEP] a search for knowledge [SEP]	No	research%1:09:00::
[CLS] Your research ... [SEP] inquire into [SEP]	No	research%2:31:00::
[CLS] Your research ... [SEP] attempt to find out in a ... [SEP]	No	research%2:32:00::

Context-Gloss Pairs with weak supervision of the target word [research]

	Label	Sense Key
[CLS] Your “research” ... [SEP] research: systematic investigation to ... [SEP]	Yes	research%1:04:00::
[CLS] Your “research” ... [SEP] research: a search for knowledge [SEP]	No	research%1:09:00::
[CLS] Your “research” ... [SEP] research: inquire into [SEP]	No	research%2:31:00::
[CLS] Your “research” ... [SEP] research: attempt to find out in a ... [SEP]	No	research%2:32:00::

Table 1: The construction methods. The sentence is taken from SemEval-2007 WSD dataset. The ellipsis “...” indicates the remainder of the sentence or the gloss.

from glosses of all possible senses (in WordNet) of the target word, thus treating WSD task as a sentence-pair classification problem. We fine-tune the pre-trained BERT model on SemCor3.0 training corpus.

Recently², we are informed by Vial et al. (2019)³ that they also use BERT and incorporate WordNet as lexical knowledge in their supervised WSD system. But our work is much different from theirs. They exploit the semantic relationships between senses such as synonymy, hypernymy and hyponymy and rely on pre-trained BERT word vectors (feature-based approach); we leverage gloss knowledge (sense definition) and use BERT through fine-tuning procedures. Out of respect, we add their results in Table 3. However, the results of their feature-based approach in the same experimental setup (single training set and single model) are not as good as our fine-tuning approach although their ensemble systems (with another training set WNGC) achieve better performance.

In particular, our contribution is two-fold:

1. We construct *context-gloss* pairs and propose three BERT-based models for WSD.
2. We fine-tune the pre-trained BERT model on SemCor3.0 training corpus, and the experimental results on several English all-words WSD benchmark datasets show that our approach outperforms the state-of-the-art systems.

2 Methodology

In this section, we describe our method in detail.

²after we submitted the final version to the conference

³their paper is available on arXiv after our first submission to the conference in May, 2019

2.1 Task Definition

In WSD, a sentence s usually consists of a series of words: $\{w_1, \dots, w_m\}$, and some of the words $\{w_{i_1}, \dots, w_{i_k}\}$ are targets $\{t_1, \dots, t_k\}$ need to be disambiguated. For each target t , its candidate senses $\{c_1, \dots, c_n\}$ come from entries of its lemma in a pre-defined sense inventory (usually WordNet). Therefore, WSD task aims to find the most suitable entry (symbolized as unique sense key) for each target in a sentence. See a sentence example in Table 1.

2.2 BERT

BERT (Devlin et al., 2018) is a new language representation model, and its architecture is a multi-layer bidirectional Transformer encoder. BERT model is pre-trained on a large corpus and two novel unsupervised prediction tasks, i.e., masked language model and next sentence prediction tasks are used in pre-training. When incorporating BERT into downstream tasks, the fine-tuning procedure is recommended. We fine-tune the pre-trained BERT model on WSD task.

BERT(Token-CLS) Since every target in a sentence needs to be disambiguated to find its exact sense, WSD task can be regarded as a token-level classification task. To incorporate BERT to WSD task, we take the final hidden state of the token corresponding to the target word (if more than one token, we average them) and add a classification layer for every target lemma, which is the same as the last layer of the Bi-LSTM model (Kågebäck and Salomonsson, 2016).

2.3 GlossBERT

BERT can explicitly model the relationship of a pair of texts, which has shown to be beneficial to many pair-wise natural language understanding tasks. In order to fully leverage gloss information, we propose GlossBERT to construct *context-gloss* pairs from all possible senses of the target word in WordNet, thus treating WSD task as a sentence-pair classification problem.

We describe our construction method with an example (See Table 1). There are four targets in this sentence, and here we take target word *research* as an example:

Context-Gloss Pairs The sentence containing target words is denoted as *context* sentence. For each target word, we extract glosses of all N possible senses (here $N = 4$) of the target word (research) in WordNet to obtain the *gloss* sentence. [CLS] and [SEP] marks are added to the *context-gloss* pairs to make it suitable for the input of BERT model. A similar idea is also used in aspect-based sentiment analysis (Sun et al., 2019).

Context-Gloss Pairs with Weak Supervision Based on the previous construction method, we add weak supervised signals to the *context-gloss* pairs (see the highlighted part in Table 1). The signal in the *gloss* sentence aims to point out the target word, and the signal in the *context* sentence aims to emphasize the target word considering the situation that a target word may occur more than one time in the same sentence.

Therefore, each target word has N *context-gloss* pair training instances ($label \in \{yes, no\}$). When testing, we output the probability of $label = yes$ of each *context-gloss* pair and choose the sense corresponding to the highest probability as the prediction label of the target word. We experiment with three GlossBERT models:

GlossBERT(Token-CLS) We use context-gloss pairs as input. We highlight the target word by taking the final hidden state of the token corresponding to the target word (if more than one token, we average them) and add a classification layer ($label \in \{yes, no\}$).

GlossBERT(Sent-CLS) We use context-gloss pairs as input. We take the final hidden state of the first token [CLS] as the representation of the whole sequence and add a classification layer

Dataset	Total	Noun	Verb	Adj	Adv
SemCor	226036	87002	88334	31753	18947
SE2	2282	1066	517	445	254
SE3	1850	900	588	350	12
SE07	455	159	296	0	0
SE13	1644	1644	0	0	0
SE15	1022	531	251	160	80

Table 2: Statistics of the different parts of speech annotations in English all-words WSD datasets.

($label \in \{yes, no\}$), which does not highlight the target word.

GlossBERT(Sent-CLS-WS) We use context-gloss pairs with weak supervision as input. We take the final hidden state of the first token [CLS] and add a classification layer ($label \in \{yes, no\}$), which weekly highlight the target word by the weak supervision.

3 Experiments

3.1 Datasets

The statistics of the WSD datasets are shown in Table 2.

Training Dataset Following previous work (Luo et al., 2018a,b; Raganato et al., 2017a,b; Iacobacci et al., 2016; Zhong and Ng, 2010), we choose SemCor3.0 as training corpus, which is the largest corpus manually annotated with WordNet sense for WSD.

Evaluation Datasets We evaluate our method on several English all-words WSD datasets. For a fair comparison, we use the benchmark datasets proposed by Raganato et al. (2017b) which include five standard all-words fine-grained WSD datasets from the Senseval and SemEval competitions: Senseval-2 (SE2), Senseval-3 (SE3), SemEval-2007 (SE07), SemEval-2013 (SE13) and SemEval-2015 (SE15). Following Luo et al. (2018a), Luo et al. (2018b) and Raganato et al. (2017a), we choose SE07, the smallest among these test sets, as the development set.

WordNet Since Raganato et al. (2017b) map all the sense annotations in these datasets from their original versions to WordNet 3.0, we extract word sense glosses from WordNet 3.0.

	Dev	Test Datasets				Concatenation of Test Datasets				
System	SE07	SE2	SE3	SE13	SE15	Noun	Verb	Adj	Adv	All
MFS baseline	54.5	65.6	66.0	63.8	67.1	67.7	49.8	73.1	80.5	65.5
Lesk _{ext+emb}	56.7	63.0	63.7	66.2	64.6	70.0	51.1	51.7	80.6	64.2
Babelfy	51.6	67.0	63.5	66.4	70.3	68.9	50.7	73.2	79.8	66.4
IMS	61.3	70.9	69.3	65.3	69.5	70.5	55.8	75.6	82.9	68.9
IMS+emb	62.6	72.2	70.4	65.9	71.5	71.9	56.6	75.9	84.7	70.1
Bi-LSTM	-	71.1	68.4	64.8	68.3	69.5	55.9	76.2	82.4	68.4
Bi-LSTM+att.+LEX+POS	64.8	72.0	69.1	66.9	71.5	71.5	57.5	75.0	83.8	69.9
GAS _{ext} (Linear)	-	72.4	70.1	67.1	72.1	71.9	58.1	76.4	84.7	70.4
GAS _{ext} (Concatenation)	-	72.2	70.5	67.2	72.6	72.2	57.7	76.6	85.0	70.6
CAN ^s	-	72.2	70.2	69.1	72.2	73.5	56.5	76.6	80.3	70.9
HCAN	-	72.8	70.3	68.5	72.8	72.7	58.2	77.4	84.1	71.1
SemCor, hypernoms (single)	-	-	-	-	-	-	-	-	-	75.6
SemCor, hypernoms (ensemble)†	69.5	77.5	77.4	76.0	78.3	79.6	65.9	79.5	85.5	76.7
SemCor+WNGC, hypernoms (single)‡	-	-	-	-	-	-	-	-	-	77.1
SemCor+WNGC, hypernoms (ensemble)† ‡	73.4	79.7	77.8	78.7	82.6	81.4	68.7	83.7	85.5	79.0
BERT(Token-CLS)	61.1	69.7	69.4	65.8	69.5	70.5	57.1	71.6	83.5	68.6
GlossBERT(Sent-CLS)	69.2	76.5	73.4	75.1	79.5	78.3	64.8	77.6	83.8	75.8
GlossBERT(Token-CLS)	71.9	77.0	75.4	74.6	79.3	78.3	66.5	78.6	84.4	76.3
GlossBERT(Sent-CLS-WS)	72.5	77.7	75.2	76.1	80.4	79.3	66.9	78.2	86.4	77.0

Table 3: F1-score (%) for fine-grained English all-words WSD on the test sets in the framework of Raganato et al. (2017b) (including the development set SE07). The six blocks list the MFS baseline, two knowledge-based systems, two traditional *word expert* supervised systems, six recent neural-based systems, one BERT feature-based system and our systems, respectively. Results in first three blocks come from Raganato et al. (2017b), and others from the corresponding papers. [†] values are ensemble systems and [‡] values are models trained on both SemCor and WNGC. **Bold** font indicates best single model system trained on SemCor, i.e. excludes ^{† ‡} values since it is meaningless to compare ensemble systems and models trained on two training sets with our single model trained on SemCor training set only.

3.2 Settings

We use the pre-trained uncased BERT_{BASE} model⁴ for fine-tuning, because we find that BERT_{LARGE} model performs slightly worse than BERT_{BASE} in this task. The number of Transformer blocks is 12, the number of the hidden layer is 768, the number of self-attention heads is 12, and the total number of parameters of the pre-trained model is 110M. When fine-tuning, we use the development set (SE07) to find the optimal settings for our experiments. We keep the dropout probability at 0.1, set the number of epochs to 4. The initial learning rate is 2e-5, and the batch size is 64.

3.3 Results

Table 3 shows the performance of our method on the English all-words WSD benchmark datasets. We compare our approach with previous methods.

The first block shows the MFS baseline, which selects the most frequent sense in the training corpus for each target word.

The second block shows two knowledge-based systems. Lesk_{ext+emb} (Basile et al., 2014) is a variant of Lesk algorithm (Lesk, 1986) by calcu-

lating the gloss-context overlap of the target word. Babelfy (Moro et al., 2014) is a unified graph-based approach which exploits the semantic network structure from BabelNet.

The third block shows two *word expert* traditional supervised systems. IMS (Zhong and Ng, 2010) is a flexible framework which trains SVM classifiers and uses local features. And IMS_{+emb} (Iacobacci et al., 2016) is the best configuration of the IMS framework, which also integrates word embeddings as features.

The fourth block shows several recent neural-based methods. Bi-LSTM (Kågebäck and Salomonsson, 2016) is a baseline for neural models. Bi-LSTM_{+att.+LEX+POS} (Raganato et al., 2017a) is a multi-task learning framework for WSD, POS tagging, and LEX with self-attention mechanism, which converts WSD to a sequence learning task. GAS_{ext} (Luo et al., 2018b) is a variant of GAS which is a gloss-augmented variant of the memory network by extending gloss knowledge. CAN^s and HCAN (Luo et al., 2018a) are sentence-level and hierarchical co-attention neural network models which leverage gloss knowledge.

The fifth block are feature-based BERT models (Vial et al., 2019) which exploit the semantic relationships between senses such as synonymy,

⁴https://storage.googleapis.com/bert_models/2018_10_18/uncased_L-12_H-768_A-12.zip

hypernymy and hyponymy, and use pre-trained BERT embeddings and transformer encoder layers. It is worth noting that our fine-tuning method is superior to their feature-based method under the same experimental setup (single model + SemCor training set).

In the last block, we report the performance of our method. BERT(Token-CLS) is our baseline, which does not incorporate gloss information, and it performs slightly worse than previous traditional supervised methods and recent neural-based methods. It proves that directly using BERT cannot obtain performance growth. The other three methods outperform other models by a substantial margin, which proves that the improvements come from leveraging BERT to better exploit gloss information. It is worth noting that our method achieves significant improvements in **SE07** and **Verb** over previous methods, which have the highest ambiguity level among all datasets and all POS tags respectively according to Raganato et al. (2017b).

Moreover, GlossBERT(Token-CLS) performs better than GlossBERT(Sent-CLS), which proves that highlighting the target word in the sentence is important. However, the weakly highlighting method GlossBERT(Sent-CLS-WS) performs best in most circumstances, which may result from its combination of the advantages of the other two methods.

3.4 Discussion

There are two main reasons for the great improvements of our experimental results. First, we construct *context-gloss* pairs and convert WSD problem to a sentence-pair classification task which is similar to NLI tasks and train only one classifier, which is equivalent to expanding the corpus. Second, we leverage BERT (Devlin et al., 2018) to better exploit the gloss information. BERT model shows its advantage in dealing with sentence-pair classification tasks by its amazing improvement on QA and NLI tasks. This advantage comes from both of its two novel unsupervised prediction tasks.

Compared with traditional *word expert* supervised methods, our GlossBERT shows its effectiveness to alleviate the effort of feature engineering and does not require training a dedicated classifier for every target lemma. Compared with recent neural-based methods, our solution is more intuitive and can make better use of gloss knowl-

edge. Besides, our approach demonstrates that when we fine-tune BERT on a downstream task, converting it into a sentence-pair classification task may be a good choice.

4 Conclusion

In this paper, we seek to better leverage gloss knowledge in a supervised neural WSD system. We propose a new solution to WSD by constructing *context-gloss* pairs and then converting WSD to a sentence-pair classification task. We fine-tune the pre-trained BERT model on SemCor3.0 training corpus, and the experimental results on several English all-words WSD benchmark datasets show that our approach outperforms the state-of-the-art systems.

Acknowledgments

We would like to thank the anonymous reviewers for their valuable comments. The research work is supported by National Natural Science Foundation of China (No. 61751201 and 61672162), Shanghai Municipal Science and Technology Commission (16JC1420401 and 17JC1404100), Shanghai Municipal Science and Technology Major Project (No.2018SHZDZX01) and ZJLab.

References

- Eneko Agirre, Oier López de Lacalle, and Aitor Soroa. 2014. Random walks for knowledge-based word sense disambiguation. *Computational Linguistics*, 40(1):57–84.
- Satanjeev Banerjee and Ted Pedersen. 2002. An adapted lesk algorithm for word sense disambiguation using wordnet. In *International conference on intelligent text processing and computational linguistics*, pages 136–145. Springer.
- Pierpaolo Basile, Annalina Caputo, and Giovanni Semeraro. 2014. An enhanced lesk word sense disambiguation algorithm through a distributional semantic model. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1591–1600.
- 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*.
- Ignacio Iacobacci, Mohammad Taher Pilehvar, and Roberto Navigli. 2016. Embeddings for word sense disambiguation: An evaluation study. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 897–907.

- Mikael Kågebäck and Hans Salomonsson. 2016. Word sense disambiguation using a bidirectional lstm. *arXiv preprint arXiv:1606.03568*.
- Michael Lesk. 1986. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *Proceedings of the 5th annual international conference on Systems documentation*, pages 24–26. Citeseer.
- Fuli Luo, Tianyu Liu, Zexue He, Qiaolin Xia, Zhifang Sui, and Baobao Chang. 2018a. Leveraging gloss knowledge in neural word sense disambiguation by hierarchical co-attention. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1402–1411.
- Fuli Luo, Tianyu Liu, Qiaolin Xia, Baobao Chang, and Zhifang Sui. 2018b. Incorporating glosses into neural word sense disambiguation. *arXiv preprint arXiv:1805.08028*.
- George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.
- Andrea Moro, Alessandro Raganato, and Roberto Navigli. 2014. Entity linking meets word sense disambiguation: a unified approach. *Transactions of the Association for Computational Linguistics*, 2:231–244.
- Roberto Navigli. 2009. Word sense disambiguation: A survey. *ACM computing surveys (CSUR)*, 41(2):10.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237.
- Alessandro Raganato, Claudio Delli Bovi, and Roberto Navigli. 2017a. Neural sequence learning models for word sense disambiguation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1156–1167.
- Alessandro Raganato, Jose Camacho-Collados, and Roberto Navigli. 2017b. Word sense disambiguation: A unified evaluation framework and empirical comparison. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 99–110.
- Hui Shen, Razvan Bunescu, and Rada Mihalcea. 2013. Coarse to fine grained sense disambiguation in wikipedia. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, volume 1, pages 22–31.
- Chi Sun, Luyao Huang, and Xipeng Qiu. 2019. Utilizing bert for aspect-based sentiment analysis via constructing auxiliary sentence. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 380–385.
- Loïc Vial, Benjamin Lecouteux, and Didier Schwab. 2019. Sense vocabulary compression through the semantic knowledge of wordnet for neural word sense disambiguation. *arXiv preprint arXiv:1905.05677*.
- Zhi Zhong and Hwee Tou Ng. 2010. It makes sense: A wide-coverage word sense disambiguation system for free text. In *Proceedings of the ACL 2010 system demonstrations*, pages 78–83.