

Decomposing Word Embedding with the Capsule Network

Xin Liu^{a,1}, Qingcai Chen^{a,*}, Yan Liu^b, Joanna Siebert^a, Baotian Hu^a,
Xiangping Wu^a, Buzhou Tang^a

^a*Harbin Institute of Technology, Shenzhen*

^b*Department of Computing, The Hong Kong Polytechnic University, Hong Kong*

Abstract

Word sense disambiguation tries to learn the appropriate sense of an ambiguous word in a given context. The existing pre-trained language methods and the methods based on multi-embeddings of word did not explore the power of the unsupervised word embedding sufficiently.

In this paper, we discuss a capsule network-based approach, taking advantage of capsules potential for recognizing highly overlapping features and dealing with segmentation. We propose a **Capsule** network-based method to **Decompose** the unsupervised word **E**mbedding of an ambiguous word into context specific **S**ense embedding, called CapsDecE2S. In this approach, the unsupervised ambiguous embedding is fed into capsule network to produce its multiple morpheme-like vectors, which are defined as the basic semantic language units of meaning. With attention operations, CapsDecE2S integrates the word context to reconstruct the multiple morpheme-like vectors into the context-specific sense embedding. To train CapsDecE2S, we propose a sense matching training method. In this method, we convert the sense learning into a binary classification that explicitly learns the relation between senses by the label of matching and non-matching. The CapsDecE2S was experimentally evaluated on two sense learning tasks, i.e., word in context and word sense disambiguation. Results on two public corpora Word-in-Context and English all-words Word Sense Disambiguation show that, the CapsDecE2S model achieves the new state-of-the-art for the word in context and word sense disambiguation tasks.

Keywords: Word sense learning, Capsule network, Word-in-Context, English all-words Word Sense Disambiguation

*Corresponding author

Email addresses: hit.liuxin@gmail.com (Xin Liu), qingcai.chen@gmail.com (Qingcai Chen), csyliu@comp.polyu.edu.hk (Yan Liu), joannasiebert@yahoo.com (Joanna Siebert), baotiannlp@gmail.com (Baotian Hu), wxpleduole@gmail.com (Xiangping Wu), tangbuzhou@gmail.com (Buzhou Tang)

¹orcid=https://orcid.org/0000-0003-2802-594X

1. Introduction

Word meanings are determined by their contexts [1]. It is a generally followed principle by unsupervised word embedding approaches [2], e.g., Word2vec [3] and GloVe [4]. In these approaches, the words with similar semantic roles are mapped into nearby points in the embedding space. They successfully capture the word semantic properties. However, there is still an important issue left, i.e., distinguishing between the multiple senses of an ambiguous word whose intended meaning is in some way unclear to the reader because of its multiple interpretations or definitions.

The ambiguous words usually present multiple senses in various contexts. Many approaches for learning the appropriate sense of an ambiguous word in a given context have been proposed. The pre-trained language models, e.g., ELMo [5], OPENAI GPT [6], and BERT [7], are popular methods for learning the word sense representation dynamically. The word embeddings in ELMo are learned functions of the internal states of a deep bidirectional language model [5]. The OPENAI GPT uses a left-to-right Transformer and is trained on the BooksCorpus (800M words) [6] while BERT uses the pre-trained language tasks to learn word representations with the transformer language model [7]. These methods output contextual embeddings that infer different representations induced by arbitrarily long contexts. They have had a major impact on driving progress on downstream tasks [8]. However, no explicit sense label is involved in these methods.

The multi-embeddings of an ambiguous word are another popular solution for sense learning. The ambiguous word is usually represented by multiple embeddings and each embedding corresponds to one of its senses. Sense graphs [9, 10], bilingual resources [11, 12, 13, 14], and semantic network [15, 16, 17, 18] are widely used for learning multiple embeddings. [9] proposed to apply graph smoothing and predictive maximum likelihood models to learn senses grounded in a specified ontology. [14] proposed to retrofit sense-specific word vectors using parallel text. [19] extracted semantically related words from WordNet and computed the sense embeddings in turn.

The multi-embedding based methods usually require well-organized knowledge base, whose scale is usually smaller than that for unsupervised word embedding learning. Then, a natural question emerges: can we learn some information for the proper word sense based on the unsupervised word embedding? For example, in " \mathcal{S}_1 : Which fruit contains more vitamin C, *apple* or strawberry" and " \mathcal{S}_2 : *Apple* is about to release iPhone X", the embedding "apple" gives higher similarities to the words related to one of its senses than others ("strawberry" and "iPhone", respectively) This phenomenon indicates that we may infer some exact sense information from the unsupervised word embedding [20].

However, the aforementioned example also shows that the unsupervised word embedding of an ambiguous word contains some redundant information from other senses. Then an urgent problem is how to extract the useful information from the unsupervised embedding that simultaneously contains other senses. For this problem, we believe that using the capsule network could be a solution.

The capsule network was first proposed by [21] for digital image recognition task, and it shows potential for recognizing highly overlapping features [21, 22]. In capsule network, a capsule is a group of neurons that represent the instantiation parameters of the entity [21], which means that a capsule is a vector. Capsules are also very good for dealing with segmentation. A lower-level capsule prefers to hand out its output to higher-level capsules and the higher-level capsules select information from the lower-level capsules to represent more complex entities [21]. At each location in the feature, there is at most one instance of the type of entity that a capsule represents [21]. In other words, each capsule represents one kind of entity which contains the unique information. The ambiguous word embedding is just like the highly overlapping feature and can take full advantage of the capsule network approach. When the embedding goes through the capsule network, the capsule at one level attends to some active capsules at the level below and ignores others. Finally, the higher-level capsules contains the unique information, and the capsules are just like the morphemes since each represents one basic semantic unit of meaning. The procedure acts like decomposing an ambiguous word embedding into multiple morpheme-like vectors. Even though the capsule network has been widely used for image recognition, we are unaware of any capsule network based approaches for word sense disambiguation. We believe that capsule network can benefit the word sense disambiguation, especially decomposing the unsupervised word embedding.

However, applying the capsule network to embedding decomposing is not a trivial task. It is challenging how to allocate the weights for combining the decomposed morpheme-like vectors. Following the general principle that word meaning is determined by its context [1], a context attention operation can be used to compute the weights. Attention is the widely used mechanism to draw global dependencies between input and output [23]. Taking the context as input and each decomposed morpheme-like vector as output, the attention weights can be easily obtained. By combining these morpheme-like vectors with different weights, we can reach different senses of the same word.

For word sense learning, some approaches learn the exact sense by a classification training with the senses as the labels [24, 25, 26]. Some other works proposed to learn the word sense by k-nearest neighbor cluster training [8, 27, 28]. However, there are large amount of labels in the first group of methods and it is hard to train a model with so many labels. On the other hand, the cluster training-based methods may face the problem of inexact senses when clustering. Besides, both methods cannot solve the sense zero-shot problem in training and test set. Inspired by the recent works [29, 30, 31], we propose a flexible training method, called word sense matching training, which converts the problem to a binary classification by judging the matching state of two senses. With this training method, though the same sense of the ambiguous word may cross different sentences, we still learn the appropriate sense representation. Following [29], the gloss of word sense can be converted into a sentence with the form of "word:gloss", which provides the gold standard for matching. Besides, since the gloss covers all the senses in WordNet, we can match almost all senses and

obstacle the problem of zero-shot. On one hand, unlike the previous classification training, the word sense matching method only outputs the matching label of two senses. On the other hand, the method uses the explicit sense as the gold standard for training and the zero-shot could be solved by matching the glosses from WordNet [30].

Our main contributions are summarized as follows. 1) We propose an embedding decomposing method with the capsule network. The capsule network decomposes the unsupervised word embedding into multiple morpheme-like vectors. 2) We propose a novel framework for merging morpheme-like vectors by contextual attention, which enables the dynamic generation of context specific sense representation. The context specific sense representation maintains the appropriate sense of an ambiguous word in a given context. 3) We propose a new training method called word sense matching for word sense learning, which integrates the gold sense label into the training procedure, and achieves the new state-of-the-art on word in context task and word sense disambiguation task, respectively.

2. Related Works

Many efforts have been made for word sense learning. In terms of resources, several methods for multiple senses representation learning automatically induce word senses from monolingual corpora [11, 32]. [11] provided a context-dependent vector representation of word meaning with the Wikipedia and Gigaword corpus.[32] uses the automatic discovery of word senses to Web search result clustering. Others extended it to using bilingual corpora [12, 13, 14, 33, 34, 35]. [12] proposed to learn sense-specific word embeddings by exploiting bilingual resources. [13] presented an extension to the Skip-gram model to learn multiple embeddings per word type. [14] proposed to retrofit sense-specific word vectors using parallel text. [33] used bilingual learning of multi-sense embeddings with discrete autoencoders. These methods focus on the statistics information extracted from text corpora and contribute to the development of the word sense learning research area greatly. However, at the same time, they ignore exploring knowledge from semantic networks [17]. As a result, the induced senses are not readily interpretable and are not easily mappable to lexical resources either [36].

By taking advantage of the information in a knowledge base, the knowledge-based approaches are proposed to exploit knowledge resources like WordNet and BabelNet [15, 16, 17, 18, 30, 37, 38]. [15] presented a WSD algorithm based on random walks over large Lexical Knowledge Bases (LKB). [16] analyzed how using a resource that integrates BabelNet might enable WSD to be solved. [18] presented two fully automatic and language-independent sense computing methods based on BabelNet and Wikipedia. [17] exploited large corpora and knowledge from the semantic networks to produce word and sense embeddings. The multi-embeddings based methods usually rely on the knowledge base when learning multiple embeddings for each sense [9, 10, 11]. The advantage of the knowledge-based systems is that they do not require sense-annotated data, and

the knowledge enhances the word sense learning ability of these methods. However, without the sense-annotated data in these methods, the performance is also limited, and each disambiguation word is treated in isolation with a weak relationship [26]).

In order to address deficiencies of the knowledge-based approaches, some works exploiting the use of the sense annotated corpus have been proposed. In works that use supervised data, a machine learning classifier is trained with a large amount of data with annotated senses [39, 40, 26]. Usually, they depend greatly on the annotated corpus, e.g. SemCor3.0 and OMSTI. It is at the expense, however, of harder training and limited flexibility [41].

As mentioned above, the monolingual and bilingual corpora, the knowledge base, and the sense annotated corpus improve the word sense learning ability of the methods. Combining these resources can further improve the performance of sense learning methods. For example, even though the BERT model [7] shows outstanding performance on word sense learning due to the language pre-training, even better performance can be reached when the knowledge base or supervised corpus is integrated into the BERT model [8, 29, 28]. [8] focused on the synset, hypernym, and lexname with full-coverage of WordNet on BERT. [29] focused on better leveraging gloss knowledge into the BERT model and used the SemCor3.0 corpus to train the model. Both methods have provided the new state-of-the-art on the word sense learning task in succession.

All the above methods have made a great contribution to the process of word sense learning. Our work is highly related to the previous ones but we propose a word sense learning method from a new perspective. In our proposed method, we learn the appropriate sense of an ambiguous word from the unsupervised word embedding with the powerful capsule network and contextual attention, and the method is trained with the sense-annotated corpus by word sense matching training. We have found that using unsupervised word embedding with the capsule network can reach the new state-of-the-art on word sense learning tasks. However, we are unaware of any such methods in word sense learning.

3. Method

Figure 1 depicts an overall diagram of the proposed CapsDecE2S method. There are two main modules in the CapsDecE2S method: the embedding decomposing and the context learning.

In CapsDecE2S, the ambiguous word in the sentence is regarded as the target word. With the embedding lookup, we can derive each word embedding in the sentence. The target word embedding is fed into the embedding decomposing module and to the context learning module. The embedding decomposing module decomposes the target word embedding with the capsule network and outputs the multiple morpheme-like vectors. At the same time, in the context learning module, each word embedding in the sentence is used to first derive the global and local context representation, and then together with the decomposed morpheme-like vectors from the decomposition module to learn the context-specific sense representation of the target word.

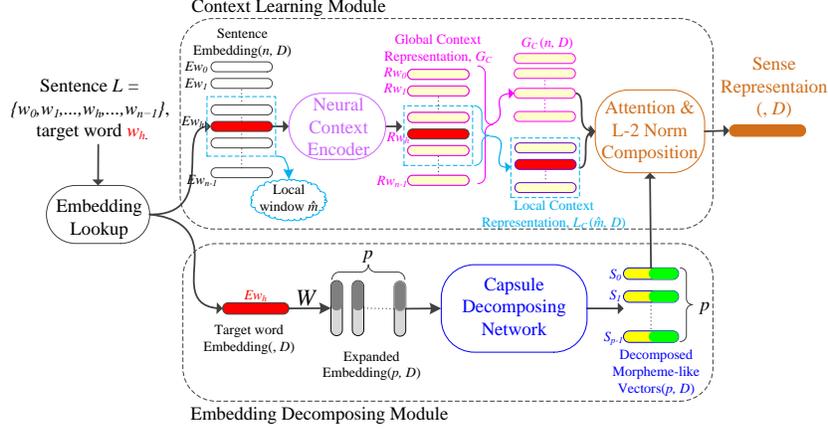


Figure 1: The embedding decomposing and context learning procedure of the CapsDecE2S model. The numbers in the bracket mean the variable dimension.

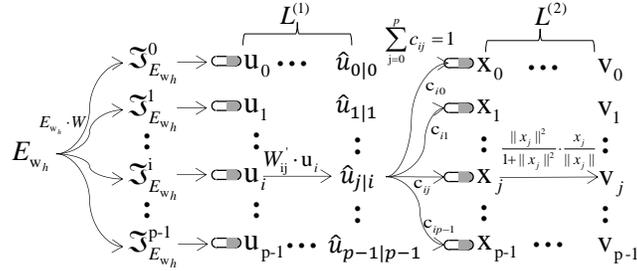


Figure 2: The decomposing calculation between capsules in the initial two layers. In the following layer, the outputs $\{v_0, v_1, \dots, v_{p-1}\}$ are used as its input, namely $u_k = v_k, k \in [0, p)$.

In Section 3.1, we describe how to decompose the unsupervised embedding with the capsule network in detail and in section 3.2, we describe the context learning procedure. In Section 3.3, we introduce the proposed word sense matching training on how to train the CapsDecE2S method.

3.1. The Embedding Decomposing with the Capsule Network

In the embedding decomposing module shown in Figure 1, the target word embedding is first expanded by parameters W . Next, the expanded embeddings are input into the capsule decomposing network. Then, the capsule decomposing network outputs the decomposed morpheme-like vectors.

Figure 2 depicts the decomposing calculation of two initial layers in the capsule decomposing network. $L = \{w_0, w_1, \dots, w_h, \dots, w_{n-1}\}$ denotes a sentence, and the target word w_h embedding is E_{w_h} . First, E_{w_h} is expanded with parameter W as $\mathfrak{S}_{E_{w_h}}^i = E_{w_h} \cdot W_i, i \in [0, p)$, where p is the maximum number of the decomposed vectors. In the first layer, the input is $\mathfrak{S}_{E_{w_h}}^i$, and each $\mathfrak{S}_{E_{w_h}}^i$

corresponds to one capsule. For a capsule i in the layer $L^{(1)}$ (abbr. $i_{L^{(1)}}$), we have $u_i = \mathfrak{S}_{E_{w_h}}^i$. A weight matrix W'_{ij} is used for building connections with the capsule j in the layer $L^{(2)}$ (abbr. $j_{L^{(2)}}$), and a prediction vector $\hat{u}_{j|i}$ is produced. Next, the total input x_j to the capsule $j_{L^{(2)}}$ is a weighted sum over all $\hat{u}_{j|i}$ from the capsules in the layer $L^{(1)}$.

$$x_j = \sum_i c_{ij} \hat{u}_{j|i}, \quad \hat{u}_{j|i} = W'_{ij} u_i, \quad (1)$$

where c_{ij} is the coupling coefficient from capsule $i_{L^{(1)}}$ to $j_{L^{(2)}}$. The coupling coefficients sum to 1 between $i_{L^{(1)}}$ and all capsules in $L^{(2)}$.

$$\sum_{j=0}^p c_{ij} = 1. \quad (2)$$

In the capsule $j_{L^{(2)}}$, a non-linear "squashing" function is applied to keep the length by shrinking short vectors to almost 0 and long vectors to a length slightly below 1, is shown in Equation 3.

$$v_j = \frac{\|x_j\|^2}{1 + \|x_j\|^2} \cdot \frac{x_j}{\|x_j\|}, \quad (3)$$

where v_j is the squashing output of the capsule $j_{L^{(2)}}$.

The coupling coefficient c_{ij} is updated by the iterative dynamic routing, and it is a softmax result based on the logic b_{ij} .

$$c_{ij} = \frac{\exp(b_{ij})}{\sum_{k=0}^p \exp(b_{ik})}, \quad (4)$$

we follow the processing by [21]. Initially, b_{ij} equals to 0 and is updated as

$$b_{ij} = b_{ij} + v_j \cdot \hat{u}_{j|i}, \quad (5)$$

which aims to measure the agreement between the output v_j of $j_{L^{(2)}}$ and the prediction $\hat{u}_{j|i}$ of $i_{L^{(1)}}$.

In the following layers, the network repeats the same calculation. The output v is passed into the capsules in the next layer and goes through the weight matrix, the weighted sum and the non-linear squashing function. With K layer iterations, we take the outputs of layer K as the decomposed vectors $\{S_0, S_1, \dots, S_{p-1}\}$, where $S_j = v_{j_{L^{(K)}}}$.

3.2. The Context Learning Module

In the context learning module shown in Figure 1, the word embeddings in the sentence are input into the neural context encoder. The output of the neural context encoder is used as the global context representation, and the nearby word representations of the target word are used as the local context representation. Finally, the decomposed morpheme-like vectors, the global context representation, and the local context representation are input into the "Attention&L2 Norm Composition", which then outputs the sense representation.

Here, we introduce the calculations in the context learning module in detail. We take the neural language model (NLM) as the neural context encoder to learn the context information.

First, the words in a sentence L are converted into $\{E_{w_0}, \dots, E_{w_{n-1}}\}$ by embedding lookup, and then passed into the neural context encoder seriatim.

Second, the hidden states of the neural units in the last layer of the neural context encoder are selected as the context representation. Here, all the word representations are regarded as the global context representation G_c , namely $G_c = \{R_{w_0}, \dots, R_{w_{n-1}}\}$.

Third, we extract the nearby word representations of the target word w_h as the local context representation L_c with a window size \hat{m} , namely $L_c = \{R_{w_e}, \dots, R_{w_h}, \dots, R_{w_z}\}$, where $e = \min(h - \hat{m}/2, 0)$, $z = \max(h + \hat{m}/2, n)$.

Forth, the "Attention&L-2 Norm Composition" component first calculates the global context attention weight and the local context attention weight on the decomposed morpheme-like vectors $\{S_0, S_1, \dots, S_{p-1}\}$ based on the global context representation G_c and the local context representation L_c . The global attention weight $a^G = \{a_0^G, \dots, a_k^G, \dots, a_{p-1}^G\}$ is calculated as

$$a_k^G = \frac{\exp(\hat{c}_k)}{\sum_{j=0}^n \exp(\hat{c}_{k_j})}, \quad \hat{c}_{k_j} = S_k \cdot G_{c_j}, \quad (6)$$

where S_k is one decomposed morpheme-like vector in $\{S_0, S_1, \dots, S_{p-1}\}$ and G_{c_j} is the j -th representation in the global context representation G_c . The local attention weight $a^L = \{a_0^L, \dots, a_k^L, \dots, a_{p-1}^L\}$ is also calculated in a similar way as

$$a_k^L = \frac{\exp(\hat{c}'_k)}{\sum_{j=0}^n \exp(\hat{c}'_{k_j})}, \quad \hat{c}'_{k_j} = S_k \cdot L_{c_j}, \quad (7)$$

where L_{c_j} is the j -th representation in the local context representation L_c .

Next, we apply the attention weights a^G and a^L to its global and local context representations respectively, and obtain the context-specific vectors \mathcal{S}^* as

$$\mathcal{S}_k^* = S_k + \sum_{i=0}^n a_{k_i}^G \cdot G_{c_i} + \sum_{i=0}^{z-e+1} a_{k_i}^L \cdot L_{c_i}. \quad (8)$$

Finally, we use the L-2 norm of each \mathcal{S}^* to represent the weight $\hat{b} = \{\hat{b}_0, \dots, \hat{b}_k, \dots, \hat{b}_{p-1}\}$ in composing the final sense representation Q_S in its context.

$$Q_S = \sum_{k=0}^p \hat{b}_k \cdot \mathcal{S}_k^*, \quad \hat{b} = \frac{\exp(\hat{b})}{\sum_{k=0}^p \exp(\hat{b}_k)}, \quad \hat{b}_k = \|\mathcal{S}_k^*\|^2. \quad (9)$$

3.3. Word Sense Matching Training

Figure 3 depicts the word sense matching training process to learn the context-specific sense representation. Two sentences that contain the target word are input the CapsDecE2S model respectively, and CapsDecE2S outputs the sense representation for each target word, which is the output of the

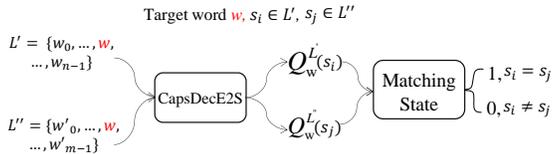


Figure 3: The word sense matching training process to predict the matching state of the word senses for w in L' and L'' , respectively. The candidate senses of word w is $\{s_0, s_1, \dots, s_m\}$ and its true senses in L' and L'' are s_i and s_j .

”Attention&L-2 Norm Composition” in context learning module. The procedure is the same as in Figure 1. Two sense representations are used to train the model by predicting their matching state. More details on the word sense matching training are given below.

In Figure 3, the sentences, e.g. L' and L'' , usually consist of a series of words. The word w is the target ambiguous word that we need to disambiguate. For w , its candidate sense set is $Set(w) = \{s_0, s_1, \dots, s_m\}$, which are defined in WordNet. The true word senses s_i and s_j in both sentences satisfy $s_i \in L'$ and $s_j \in L''$.

In word sense matching training, we define the matching and non-matching state M as:

$$M = \begin{cases} 1, & \text{if } s_i = s_j \\ 0, & \text{if } s_i \neq s_j \end{cases} \quad (10)$$

The matching prediction, based on the sense representations Q_S in L' and L'' , is used for the cross-entropy loss. First, the learned word sense representations in L' and L'' are represented as $Q_w^{L'}$ and $Q_w^{L''}$, respectively. Second, we concatenate $Q_w^{L'}$ and $Q_w^{L''}$, and pass it into the softmax classifier. A binary probability \hat{y} is predicted. Finally, based on the matching state M , the gold label is converted into binary y , and the cross-entropy loss is

$$\mathcal{L} = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})] \quad (11)$$

4. Experiments and Results

4.1. Datasets and Setup

We evaluated the proposed CapsDecE2S model on [31]’s Word-in-Context (WiC) dataset with the metric of accuracy and [42]’s cross-domain English all-words Word Sense Disambiguation (WSD) datasets (Senseval-2 (SE2) [43], Senseval-3 task1 (SE3) [44], SemEval-07 task 17 (SE07) [45], SemEval-13 task12 (SE13) [46], and SemEval-15 task13 (SE15) [16]) with the metric of F1 score.

The WiC dataset [31] is a new benchmark for the evaluation of context-sensitive word embeddings (Train:5.4K, Dev:0.63K, Test:1.4K), the target words in the dataset are nouns and verbs. On WiC, we used the default training data

and tested online ². We compare our results to the methods cited from the relevant papers [47, 8, 31].

On WSD datasets, following most researches [29, 8, 25, 42], we used the SemCor3.0 corpus as the training set and SE07 as dev set. The SemCor3.0 corpus is the largest manually annotated corpus with WordNet sense for WSD, widely used by [48, 24, 25, 42, 26, 49]. The published methods [48, 24, 26] and SOTA [8] are used for comparison.

Hyper-parameters in the experiments are as follows: the number p of the decomposed vectors (10); The capsule network layers (3); The routing iterations in capsule network (3); Unless specified otherwise, the pre-trained uncased BERT base and large are used as the neural context encoder, respectively; Especially, in WordPiece tokenization, we use the head of the sub-token as the target word. As for the others, we follow the default settings.

4.2. Baselines and Comparison Methods

The CapsDecE2S model was competed among its different versions which are implemented in this paper, along with other published methods. These published methods include the supervised models for classification, e.g., GlossAug.S. (GAS) [25], Hier.Co-Att.(HCAN) [24] and EWISE_{ConvE} [30], the supervised models for learning contextual representations of senses, e.g., Context2vec [50], ELMo [5], BERT [7] and GLU(1sent+1sur) [51], and the aggregation models with knowledge base, e.g., LMMS₂₃₄₈ [8], SemCor + hypernyms [37], and GlossBERT [29]. All these methods are under the same evaluation protocol as [31], and [42] on both tasks. The methods that used the additional manually-annotated corpus except SemCor3.0 for training and the ensemble versions are not included, e.g. partial methods in [37] and [28], but not limited to those that use the knowledge base and unsupervised resources. The proposed CapsDecE2S method and its different versions that are implemented in this paper are described as below.

- *CapsDecE2S_{base}*: This is our proposed model that all components are introduced in Section "Method", and the neural context encoder that here we used for context learning is the BERT base. The parameters for BERT base are provided by the following official link³.
- *CapsDecE2S_{large}*: The main difference from CapsDec- *E2S_{base}* is that in this model the neural context encoder for context learning is the BERT large. The parameters for BERT large are provided by the following official link⁴.
- *CapsDecE2S_{base} + LMMS*: LMMS is the sense-level embeddings with full-coverage of WordNet [8] and can be obtained from the link ⁵. In LMMS, each sense is allocated an embedding. In this model, the LMMS embedding is concatenated with the corresponding target word sense representation, $Q_{wL'}$ or

²<https://competitions.codalab.org/competitions/20010>

³https://storage.googleapis.com/bert_models/2018_10_18/uncased_L-12_H-768_A-12.zip

⁴https://storage.googleapis.com/bert_models/2018_10_18/uncased_L-24_H-1024_A-16.zip

⁵<https://github.com/danlou/lmms>

	Method	Accuracy
S-Lv	<u>1.BoW</u> [†]	58.7 [†]
	<u>2.Sentence LSTM</u> [†]	54.0 [†]
Multi-E	<u>3.DeConf</u> [†] [19]	58.7 [†]
	<u>4.SW2V</u> [†] [17]	58.1 [†]
	<u>5.JBT</u> [†] [10]	54.7 [†]
Contextual	<u>6.Context2vec</u> [†] [50]	59.3 [†]
	<u>7.ELMo</u> _{3/1} [†] [5]	58.0/57.0 [†]
	<u>8.BERT</u> _{base/large} [†] [7]	65.4/65.5 [†]
	<u>9.TextCNN</u> BERT [47]	68.6
	<u>10.LMMS</u> ₂₃₄₈ [*] [8]	69.1 [*]
	11.CapsDecE2S _{base}	70.6

Table 1: Comparisons of accuracy (%) on the WiC dataset. †- cited from WiC paper [31] with the best results and others from the corresponding papers. * - the authors only reported the dev result. The methods with an underline, wavy-line, and dash-line correspond to the sentence-level baselines (S-Lv), multi-embedding based models (Multi-E), and contextualized word based models (Contextual), respectively.

$Q_{wL''}$ in figure 3, in both sentences before they are passed into the softmax classification. First, we compute the similarity between the corresponding sentence and the LMMS sense embedding under the same lemmas, and then the most similar LMMS sense embedding is selected for concatenation. The other settings are the same as the *CapsDecE2S_{base}*.

- *CapsDecE2S_{large}*+LMMS: The difference from the *Caps-DecE2S_{base}*+LMMS is that in this model the neural context encoder for context learning is the BERT large.
- $BERT_{base}^{wsm}$: The BERT base model is fine-tuned with the word sense matching training (wsm), which the sentence pair goes through the BERT base model. Then we extract the target word representations in both sentences and concatenate them as input to the softmax classification for matching training.
- $BERT_{large}^{wsm}$: The model is trained in the same way as the $BERT_{base}^{wsm}$ model, but it is BERT large model.

4.3. Results on Word in Context

Table 1 lists the accuracies of the sentence-level baselines, multi-embedding based models, contextualized word based models, and our *CapsDecE2S_{base}* model. Since there are formal training corpus with 5.4K pairs on WiC task, we only report the *CapsDecE2S_{base}* model without using any other resources. *CapsDecE2S_{base}* outperforms all the other methods by a large margin with an accuracy of 70.6% (the *CapsDecE2S_{large}* model performs quite close to this score in this task, so we do not report it.).

The sentence-level models are widely used for sentence encoding, but they show poor performance. The main reason may be that the WiC dataset is too

	Method	SE07	SE2	SE3	SE13	SE15	ALL
	1.MFS baseline	54.5	65.6	66.0	63.8	67.1	64.8
<i>Sup</i>	2.IMS+emb [48]	62.6	72.2	70.4	65.9	71.5	69.6
	3.Seq2Seq _{multi-tasks} [26]	63.1	70.1	68.5	66.5	69.2	68.6
	4.Bi-LSTM _{multi-tasks} [26]	64.8	72.0	69.1	66.9	71.5	69.9
	5.GlossAug.S.(GAS) [25]	-	72.2	70.5	67.2	72.6	70.6
	6.Hier.Co-Att.(HCAN) [24]	-	72.8	70.3	68.5	72.8	71.1
	7.EWISE _{ConvE} [30]	67.3	73.8	71.1	69.4	74.5	71.8
	<i>Sup_c</i>	8.BERT _{Token-CLS} [29]	61.1	69.7	69.4	65.8	69.5
9.Context2Vec [50]		61.3	71.8	69.1	65.6	71.9	69.0
10.GLU(1sent+1sur) [51]		68.1	75.5	73.6	71.1	76.2	74.1
<i>Sup_c + KB</i>	11.ELMo k-NN _{full-cover.} [8]	57.1	71.5	67.5	65.3	69.6	67.9
	12.BERT k-NN _{full-cover.} [8]	66.2	76.3	73.2	71.7	74.1	73.5
	13.LMMS ₂₃₄₈ [8]	68.1	76.3	75.6	75.1	77.0	75.4
	14.SemCor+hypernyms [37]	69.5	77.5	77.4	76.0	78.3	76.7
	15.GlossBERT [29]	72.5	77.7	75.2	76.1	80.4	77.0
<i>Ours</i>	16.BERT _{base} ^{wsm}	59.2	73.3	72.2	65.3	72.6	69.8
	17.BERT _{large} ^{wsm}	59.6	73.9	72.4	65.3	73.1	70.6
	18.CapsDecE2S _{base}	67.0	77.4	76.2	75.9	77.3	76.1
	19.CapsDecE2S _{large}	68.7	78.9	77.4	75.6	77.1	76.9
	20.CapsDecE2S _{base} +LMMS	73.5	78.4	79.4	76.5	78.6	77.9
	21.CapsDecE2S _{large} +LMMS	73.8	78.8	80.7	76.6	79.4	78.6

Table 2: Comparison in terms of F1-score(%) on the English all-words WSD test sets of [42] under the manually-annotated training corpus SemCor3.0. Methods are grouped by the types of, supervised models for classification (*Sup*), supervised models for learning contextual representations of senses (*Sup_c*), the aggregation models of *Sup_c* and knowledge base (*Sup_c + KB*), and our implementations [28]. For each method, we list the version with best results.

small, and the scale limits such methods without any pre-training. The multi-embedding based models make use of external lexical resources, which helps to learn more accurate sense. The contextualized word based models benefit from the large-scale language pre-training, and thus show better performance than methods foregoing. Especially, the method LMMS₂₃₄₈ based on BERT large uses additional information from Wordnet, the dictionary embedding, and the fastText embedding. The main difference between CapsDecE2S_{base} and BERT-based models, e.g. TextCNNBERT and LMMS₂₃₄₈ (BERT), is the capsule decomposing module. CapsDecE2S_{base} contributes to the WiC dataset with more than 1.5% absolute improvement in the accuracy. The result indicates the information contained in the unsupervised word embedding has a potential under limited training data and the learning ability of the capsule network.

4.4. Results on English all-words WSD

Table 2 lists the comparison results in terms of F1-score on the English all-words WSD test sets between different types of methods published in recent years. As we can see, the CapsDecE2S based models show promising results on each test set. The CapsDecE2S_{large}+LMMS reaches the best results on nearly all test sets when compared to other approaches, and outperforms the state-of-the-art (Line 15) with 1.1% improvement on "SE2", 5.5% on "SE3", 0.5% on "SE13", and 1.6% on "ALL".

Furthermore, most "Sup_c" based models (Line 8-15) outperform the "Sup" models (Line 2-7), which proves the ability of the contextual representation of senses. However, the single BERT model shows poor performance (Line 8,16,17) no matter it is trained by context-gloss training [29], k-NN training [8], or word sense matching training. The GLU (1sent+1sur) model explores multiple strategies of integrating BERT contextualized word representations and reaches better results [51]. When the information of the knowledge base, e.g., WordNet, is integrated into the BERT model (Line 13-15), all these models show outstanding performances. For example, the GlossBERT model leverages the gloss knowledge from WordNet in a supervised neural WSD system [29]. The LMMS₂₃₄₈ model focuses on creating sense-level embeddings with full-coverage of WordNet [8]. The SemCor+hypernyms model exploits the semantic relationships, e.g., synonymy, hypernymy and hyponymy, to compress the sense vocabulary of WordNet [37].

The CapsDecE2S model is a supervised model of learning contextual representations of senses, and it explores the information from the unsupervised word embedding by the capsule network. Obviously, in this way, the model outperforms the "Sup_c" and "Sup" models and shows competitive results with the "Sup_c + KB" models. When the knowledge base embedding is integrated into the CapsDecE2S model, named CapsDecE2S+LMMS (Line 20,21), the new state-of-the-art are reached on "SE2", "SE3". "SE13" and "ALL" sets. The result on "SE15" is also encouraging.

5. Discussion

In this section, we perform analyses and schematize several examples to quantitatively interpret some properties of the CapsDecE2S model, including the ablation study on the CapsDecE2S' components, the CapsDecE2S' sense similarity across contexts, the attention weight in context, the cases that CapsDecE2S fails to learn.

5.1. Ablation Study on the CapsDecE2S' Components

First, we perform an ablation study on the CapsDecE2S' components on the English "ALL" test set in [42] to see how each of the components affects the final F1-score in Table 3. Experiments are conducted based on CapsDecE2S_{base} by removing the global context, the local context, the capsule decomposing

Method	F1	Abs/Rel(%)
CapsDecE2S _{base}	76.1	0.0 /0.0
w/o Global Context	73.6	-2.5/-3.3
w/o Word Sense Matching	72.9	-3.2/-4.2
w/o Local Context	72.8	-3.3/-4.3
w/o Capsule Network	69.9	-6.2/-8.1

Table 3: Ablation study of CapsDecE2S components on the "ALL" test set in terms of F1 score. "-" indicates the negative improvement.

module, and the word sense matching training, respectively. The "w/o capsule network" means that the expanded embeddings are directly passed into the "Attention&L-2 Norm Composition" without capsule decomposing in Figure 1. The "w/o Word Sense Matching" means the capsDecE2S method was not trained by word sense training but by the sense classification training. The last column "Abs/Rel (%)" in Table 3 gives the absolute and relative improvement, that are commonly used to measure the change or difference of two variables, on the F1 score compared to CapsDecE2S_{base}.

From Table 3, we can see that each component plays an important role in CapsDecE2S_{base}. Without the capsule network, the model drops greatly with 6.2% deterioration, and the result is quite close to BERT_{base}^{wsm} in Table 2, which proves the capsule decomposing ability. The local context seems to be more useful than the global context with F1 score of 73.6% in "w/o Global context" and 72.8% in "w/o Local context", which may be that the texts far away from the target word matter less but exist as noise in the global context. Besides, one can also see that the method with sense classification training shows worse performance than that is with the word sense matching training.

5.2. CapsDecE2S' Sense Similarity across Contexts

Second, we conduct a sense similarity validation experiment to test the robust of the CapsDecE2S sense representation, in which we measure the similarities of the CapsDecE2S sense representations when the sense occurs in different contexts. The experiment is based on the principle that even though one sense occurs in different sentences, its representations should hold high similarity.

In the experiment, we first randomly sampled twenty senses from "ALL" set, and each sense is allocated with at most five sentences. Next, we randomly paired five sentences of each sense for three times and calculated the cosine scores between the target word representations. Finally, we averaged the three cosine scores of each sense as its similarity value. In the BERT model, the target word sense representation corresponds to the hidden embedding in the last layer.

The visualized map of the sense similarity values calculated by BERT and CapsDecE2S is shown in Figure 4. The format in the column to express the word

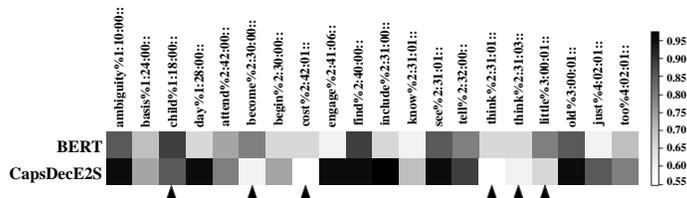


Figure 4: The visualized map of the sense similarity values from the BERT and CapsDecE2S models for sentences in "ALL" set. Column: the randomly selected senses. Color-bar: the sense similarity value scope.

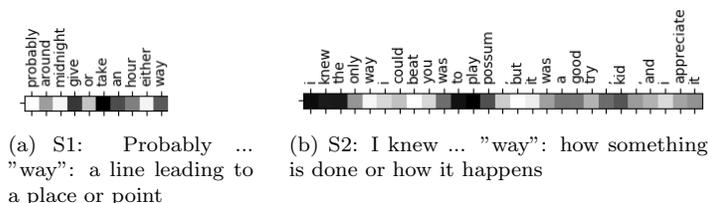


Figure 5: The attentions weights on contexts given by s entences S1 and S2 for two senses of "way", respectively. The ellipsis "..." indicates the remainder of the contexts.

sense is consistent with the definition in WordNet ⁶, and the format explanation of each field can be found here ⁷. For the contexts with more similar sense representations, their similarity value will be larger, and the color of the block in Figure 4 will be darker. From Figure 4, it is evident that apart from the 6 of 20 cases marked with a black triangle below, the remaining 16 of the blocks in the CapsDecE2S row are darker than the corresponding ones in the BERT model. Besides, we could also see that the values by the CapsDecE2S model are usually located on the upper parts in the color-bar while those by the BERT model on the lower parts. This experiment indicates that CapsDecE2S is more robust when learning the unique word sense.

5.3. Attention Weight in Context

Third, we used a noun word "way" to analyze the relationship between its sense representation and the attentive words on the sense definition.

Figure 5 schematizes the example word "way" and its two senses by analyzing the attention weights a_k^G in Equation 6 on the context. The titles of Figure 5(a) and (b) give two sense definitions from the WordNet. For either sense, we use the visualized map to present the attention weight distribution on the context when learning the context-specific sense representation. The darker block means the attention weight value on this word is larger than the lighter ones. The larger

⁶<http://wordnetweb.princeton.edu/perl/webwn>
⁷<https://wordnet.princeton.edu/documentation/senseidx5wn>

Word	Sense definition	Example Sentences
Enough	Sufficient for the purpose	But there still are not <u>enough</u> ringers to ring more than six of the eight bells.
	As much as necessary	Fortunately, these same parents do want their children to get a decent education as traditionally understood, and they have <u>enough</u> common sense to know what that demands.
Shake	Move or cause to move back and forth	This time , it just got stronger and then the building started <u>shaking</u> violently up and down.
	Move with or as if with a tremor	My back is still in knots and my hands are still <u>shaking</u> .
Plan	A series of steps to be carried out or goals to be accomplished	U.N. group drafts <u>plan</u> to reduce emissions
	The act or process of drawing up plans or layouts for some project or enterprise	There were relatively few cases reported of attempts to involve users in service <u>planning</u> but their involvement in service provision was found to be more common.

Table 4: The words for which the CapsDecE2S failed to distinguish their close senses. Example sentences are from the WSD test sets and the target words are marked with the underline.

the value is, the more the context-specific sense representation relies on this word.

From Figure 5(a) and (b), we can see that either sense relies on some words in the context, e.g. {*"give"*, *"take"*, *"an"*, *"hour"*} for sense (a) and {*"i"*, *"knew"*, *"the"*, *"only"*, *"to"*, *"play"*, *"possum"*} for (b). These words are essential in determining the unique semantic in the context, which proves that the context-specific sense representation indeed maintains the proper sense for its context.

5.4. Cases that CapsDecE2S Fails to Learn

Finally, our experiments and analyses have proven the sense learning ability of the CapsDecE2S model, but the experimental results also imply that CapsDecE2S is not omnipotent. To explore the limitation of CapsDecE2S, we collected and concluded the cases that CapsDecE2S fails to learn.

First, in all the failed cases, the top-10 failed words are the linking verbs, which include *"see"*, *"have"*, *"make"*, *"be"*, *"give"*, *"find"*, *"get"*, *"come"*, *"take"* and *"feel"*. Usually, the linking verb connects the subject with a word that gives information about the subject, such as a condition or relationship. In most cases, the linking verbs do not describe any action, instead they link the subject with the rest of the sentence. It is hard for the CapsDecE2S model to learn the linking verb’s true sense, especially since one word may occur in similar contexts. In fact, not only the CapsDecE2S model, most sense learning models are weak at these words. Second, by random sampling 10% of the failed cases, we find that except for the linking verbs, the majority are the words with quite close senses. Several typical examples are shown in Table 4. The CapsDecE2S model mistakes one sense as the other and the weeny differences between the example sentences are hard to discover.

6. Conclusion and Future Works

In this paper, we have proposed to decompose the unsupervised word embedding with the capsule network and use the context and word sense matching training to learn the sense representation. The experimental results on WiC and WSD datasets prove that the proposed CapsDecE2S method contributes to learning more accurate sense than other compared methods. These experiments show the potential of the information contained in the unsupervised word embedding and also prove the feasibility of applying the capsule network to decompose unsupervised word embedding into context-specific sense representation. Moreover, the analysis experiments show the enhanced interpretability of the capsule decomposing procedure and the context-specific sense representation.

The future works include 1) merging the CapsDecE2S representation with other sense features to improve the sense learning ability; 2) exploring the diversity of words for sense learning where the words in SemCor3.0 corpus are greatly limited by the annotation cost; 3) applying the decomposed context-specific sense representation to downstream tasks; 4) proposing solid evaluation metrics to interpret the morpheme-like vectors and context-specific sense representation.

Acknowledgments

This work is supported by Natural Science Foundation of China (Grant No. 61872113, 61573118, U1813215, 61876052), Special Foundation for Technology Research Program of Guangdong Province (Grant No. 2015B010131010), Strategic Emerging Industry Development Special Funds of Shenzhen (Grant No. JCYJ20170307150528934, JCYJ20170811153836555, JCYJ20180306172232154), Innovation Fund of Harbin Institute of Technology (Grant No. HIT.NSRIF.2017052)

References

References

- [1] J. Feng, X. Zheng, Geometric relationship between word and context representations, in: Thirty-Second AAAI Conference on Artificial Intelligence(AAAI), 2018.
- [2] R. Collobert, J. Weston, A unified architecture for natural language processing: Deep neural networks with multitask learning, in: Proceedings of the 25th international Conference on Machine Learning(ICML), 2008, pp. 160–167.
- [3] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: Advances in Neural Information Processing Systems(NeurIPS), 2013, pp. 3111–3119.

- [4] J. Pennington, R. Socher, C. Manning, Glove: Global vectors for word representation, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1532–1543. doi:10.3115/v1/D14-1162.
- [5] M. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, L. Zettlemoyer, 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), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 2227–2237. doi:10.18653/v1/N18-1202.
- [6] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, Improving language understanding by generative pre-training, in: Technical report, OpenAI., 2018.
- [7] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 4171–4186. doi:10.18653/v1/N19-1423.
- [8] D. Loureiro, A. Jorge, Language modelling makes sense: Propagating representations through WordNet for full-coverage word sense disambiguation, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Florence, Italy, 2019, pp. 5682–5691. doi:10.18653/v1/P19-1569.
- [9] S. K. Jauhar, C. Dyer, E. Hovy, Ontologically grounded multi-sense representation learning for semantic vector space models, in: Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Denver, Colorado, 2015, pp. 683–693. doi:10.3115/v1/N15-1070.
- [10] M. Pelevina, N. Arefiev, C. Biemann, A. Panchenko, Making sense of word embeddings, in: Proceedings of the 1st Workshop on Representation Learning for NLP, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 174–183. doi:10.18653/v1/W16-1620.
- [11] J. Reisinger, R. J. Mooney, Multi-prototype vector-space models of word meaning, in: Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, Los Angeles, California, 2010, pp. 109–117.

- [12] J. Guo, W. Che, H. Wang, T. Liu, Learning sense-specific word embeddings by exploiting bilingual resources, in: Proceedings of the COLING2014, 2014, pp. 497–507.
- [13] A. Neelakantan, J. Shankar, A. Passos, A. McCallum, Efficient non-parametric estimation of multiple embeddings per word in vector space, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1059–1069. doi:10.3115/v1/D14-1113.
- [14] A. Ettinger, P. Resnik, M. Carpuat, Retrofitting sense-specific word vectors using parallel text, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 1378–1383. doi:10.18653/v1/N16-1163.
- [15] E. Agirre, O. López de Lacalle, A. Soroa, Random walks for knowledge-based word sense disambiguation, Computational Linguistics 40 (1) (2014) 57–84.
- [16] A. Moro, R. Navigli, SemEval-2015 task 13: Multilingual all-words sense disambiguation and entity linking, in: Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), Association for Computational Linguistics, Denver, Colorado, 2015, pp. 288–297. doi:10.18653/v1/S15-2049.
- [17] M. Mancini, J. Camacho-Collados, I. Iacobacci, R. Navigli, Embedding words and senses together via joint knowledge-enhanced training, in: Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 100–111. doi:10.18653/v1/K17-1012.
- [18] T. Pasini, R. Navigli, Two knowledge-based methods for high-performance sense distribution learning, in: AAAI, 2018.
- [19] M. T. Pilehvar, N. Collier, De-conflated semantic representations, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1680–1690. doi:10.18653/v1/D16-1174.
- [20] E. Huang, R. Socher, C. Manning, A. Ng, Improving word representations via global context and multiple word prototypes, in: Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Jeju Island, Korea, 2012, pp. 873–882.
- [21] S. Sabour, N. Frosst, G. E. Hinton, Dynamic routing between capsules, in: NeurIPS, 2017, pp. 3856–3866.

- [22] S. Sabour, N. Frosst, G. Hinton, Matrix capsules with em routing, in: 6th international conference on learning representations, ICLR, 2018, pp. 1–15.
- [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, E. Kaiser, I. Polosukhin, Attention is all you need, in: Advances in neural information processing systems, 2017, pp. 5998–6008.
- [24] F. Luo, T. Liu, Z. He, Q. Xia, Z. Sui, B. Chang, 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, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 1402–1411. doi:10.18653/v1/D18-1170.
- [25] F. Luo, T. Liu, Q. Xia, B. Chang, Z. Sui, Incorporating glosses into neural word sense disambiguation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 2473–2482. doi:10.18653/v1/P18-1230.
- [26] A. Raganato, C. Delli Bovi, R. Navigli, Neural sequence learning models for word sense disambiguation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1156–1167. doi:10.18653/v1/D17-1120.
- [27] T. Pasini, F. Scozzafava, B. Scarlini, CluBERT: A Cluster-Based Approach for Learning Sense Distributions in Multiple Languages, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, 2020.
- [28] B. Scarlini, T. Pasini, R. Navigli, SenseBERT: Context-enhanced sense embeddings for multilingual word sense disambiguation, in: AAAI, 2020.
- [29] L. Huang, C. Sun, X. Qiu, X. Huang, GlossBERT: BERT for word sense disambiguation with gloss knowledge, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), Association for Computational Linguistics, Hong Kong, China, 2019, pp. 3507–3512. doi:10.18653/v1/D19-1355.
- [30] S. Kumar, S. Jat, K. Saxena, P. Talukdar, Zero-shot word sense disambiguation using sense definition embeddings, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Florence, Italy, 2019, pp. 5670–5681. doi:10.18653/v1/P19-1568.
- [31] M. T. Pilehvar, J. Camacho-Collados, WiC: the word-in-context dataset for evaluating context-sensitive meaning representations, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 1267–1273. doi:10.18653/v1/N19-1128.
- [32] A. Di Marco, R. Navigli, Clustering and diversifying web search results with graph-based word sense induction, *Computational Linguistics* 39 (3) (2013) 709–754. doi:10.1162/COLI_a_00148.
- [33] S. Šuster, I. Titov, G. van Noord, Bilingual learning of multi-sense embeddings with discrete autoencoders, in: *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, San Diego, California, 2016, pp. 1346–1356. doi:10.18653/v1/N16-1160.
- [34] D. Kartsaklis, M. T. Pilehvar, N. Collier, Mapping text to knowledge graph entities using multi-sense LSTMs, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 1959–1970. doi:10.18653/v1/D18-1221.
- [35] J. Camacho-Collados, M. T. Pilehvar, From word to sense embeddings: A survey on vector representations of meaning, *Journal of Artificial Intelligence Research* 63 (2018) 743–788.
- [36] A. Panchenko, E. Ruppert, S. Faralli, S. P. Ponzetto, C. Biemann, Unsupervised does not mean uninterpretable: The case for word sense induction and disambiguation, in: *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 86–98.
- [37] L. V. B. L. D. Schwab, Sense vocabulary compression through the semantic knowledge of wordnet for neural word sense disambiguation, in: *Wordnet Conference*, 2019, p. 108.
- [38] Y. Wang, M. Wang, H. Fujita, Word sense disambiguation: A comprehensive knowledge exploitation framework, *Knowledge-Based Systems* 190 (2020) 105030.
- [39] D. Kartsaklis, M. Sadrzadeh, S. Pulman, Separating disambiguation from composition in distributional semantics, in: *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, Association for Computational Linguistics, Sofia, Bulgaria, 2013, pp. 114–123.
- [40] D. Yuan, J. Richardson, R. Doherty, C. Evans, E. Altendorf, Semi-supervised word sense disambiguation with neural models, in: *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 1374–1385.

- [41] X. Liu, Q. Chen, C. Deng, H. Zeng, J. Chen, D. Li, B. Tang, LCQMC:a large-scale Chinese question matching corpus, in: Proceedings of the 27th International Conference on Computational Linguistics, Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2018, pp. 1952–1962.
- [42] A. Raganato, J. Camacho-Collados, R. Navigli, 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, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 99–110.
- [43] J. Preiss, D. Yarowsky (Eds.), Proceedings of SENSEVAL-2 Second International Workshop on Evaluating Word Sense Disambiguation Systems, Association for Computational Linguistics, Toulouse, France, 2001.
- [44] B. Snyder, M. Palmer (Eds.), Proceedings of SENSEVAL-3, the Third International Workshop on the Evaluation of Systems for the Semantic Analysis of Text, Association for Computational Linguistics, Barcelona, Spain, 2004.
- [45] S. Pradhan, E. Loper, D. Dligach, M. Palmer, SemEval-2007 task-17: English lexical sample, SRL and all words, in: Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007), Association for Computational Linguistics, Prague, Czech Republic, 2007, pp. 87–92.
- [46] R. Navigli, D. Jurgens, D. Vannella, SemEval-2013 task 12: Multilingual word sense disambiguation, in: Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013), Association for Computational Linguistics, Atlanta, Georgia, USA, 2013, pp. 222–231.
- [47] T.-Y. Chang, Y.-N. Chen, What does this word mean? explaining contextualized embeddings with natural language definition, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), Association for Computational Linguistics, Hong Kong, China, 2019, pp. 6063–6069. doi:10.18653/v1/D19-1627.
- [48] I. Iacobacci, M. T. Pilehvar, R. Navigli, Embeddings for word sense disambiguation: An evaluation study, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 897–907. doi:10.18653/v1/P16-1085.
- [49] Z. Zhong, H. T. Ng, It makes sense: A wide-coverage word sense disambiguation system for free text, in: Proceedings of the ACL 2010 System

Demonstrations, Association for Computational Linguistics, Uppsala, Sweden, 2010, pp. 78–83.

- [50] O. Melamud, J. Goldberger, I. Dagan, context2vec: Learning generic context embedding with bidirectional LSTM, in: Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 51–61. doi:10.18653/v1/K16-1006.
- [51] C. Hadiwinoto, H. T. Ng, W. C. Gan, Improved word sense disambiguation using pre-trained contextualized word representations, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), Association for Computational Linguistics, Hong Kong, China, 2019, pp. 5297–5306. doi:10.18653/v1/D19-1533.