

Incorporating Chinese Characters of Words for Lexical Sememe Prediction

Huiming Jin^{1*†}, Hao Zhu^{2†}, Zhiyuan Liu^{2,3‡}, Ruobing Xie⁴,
Maosong Sun^{2,3}, Fen Lin⁴, Leyu Lin⁴

¹ Shenyuan Honors College, Beihang University, Beijing, China

² Beijing National Research Center for Information Science and Technology,
State Key Laboratory of Intelligent Technology and Systems,
Department of Computer Science and Technology, Tsinghua University, Beijing, China

³ Jiangsu Collaborative Innovation Center for Language Ability,
Jiangsu Normal University, Xuzhou 221009 China

⁴ Search Product Center, WeChat Search Application Department, Tencent, China

Abstract

Sememes are minimum semantic units of concepts in human languages, such that each word sense is composed of one or multiple sememes. Words are usually manually annotated with their sememes by linguists, and form linguistic common-sense knowledge bases widely used in various NLP tasks. Recently, the lexical sememe prediction task has been introduced. It consists of automatically recommending sememes for words, which is expected to improve annotation efficiency and consistency. However, existing methods of lexical sememe prediction typically rely on the external context of words to represent the meaning, which usually fails to deal with low-frequency and out-of-vocabulary words. To address this issue for Chinese, we propose a novel framework to take advantage of both internal character information and external context information of words. We experiment on HowNet, a Chinese sememe knowledge base, and demonstrate that our framework outperforms state-of-the-art baselines by a large margin, and maintains a robust performance even for low-frequency words. ⁱ

1 Introduction

A sememe is an indivisible semantic unit for human languages defined by linguists (Bloomfield, 1926). The semantic meanings of concepts (e.g., words) can be composed by a finite number of sememes. However, the sememe set of a word is

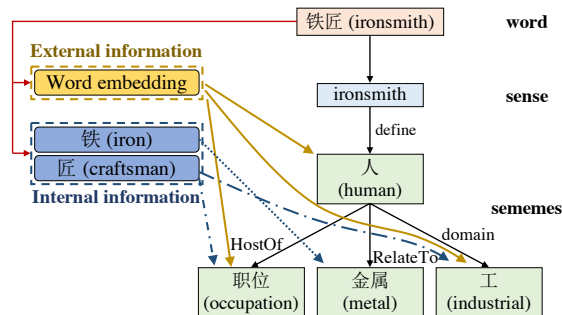


Figure 1: Sememes of the word “铁匠” (ironsmith) in HowNet, where *occupation*, *human* and *industrial* can be inferred by both external (contexts) and internal (characters) information, while *metal* is well-captured only by the internal information within the character “铁” (iron).

not explicit, which is why linguists build knowledge bases (KBs) to annotate words with sememes manually.

HowNet is a classical widely-used sememe KB (Dong and Dong, 2006). In HowNet, linguists manually define approximately 2,000 sememes, and annotate more than 100,000 common words in Chinese and English with their relevant sememes in hierarchical structures. HowNet is well developed and has a wide range of applications in many NLP tasks, such as word sense disambiguation (Duan et al., 2007), sentiment analysis (Fu et al., 2013; Huang et al., 2014) and cross-lingual word similarity (Xia et al., 2011).

Since new words and phrases are emerging every day and the semantic meanings of existing concepts keep changing, it is time-consuming and work-intensive for human experts to annotate new

* Work done while doing internship at Tsinghua University.

† Equal contribution. Huiming Jin proposed the overall idea, designed the first experiment, conducted both experiments, and wrote the paper; Hao Zhu made suggestions on ensembling, proposed the second experiment, and spent a lot of time on proofreading the paper and making revisions. All authors helped shape the research, analysis and manuscript.

‡ Corresponding author: Z. Liu (liuzy@tsinghua.edu.cn)

ⁱ Code is available at <https://github.com/thunlp/Character-enhanced-Sememe-Prediction>

concepts and maintain consistency for large-scale sememe KBs. To address this issue, Xie et al. (2017) propose an automatic sememe prediction framework to assist linguist annotation. They assumed that words which have similar semantic meanings are likely to share similar sememes. Thus, they propose to represent word meanings as embeddings (Pennington et al., 2014; Mikolov et al., 2013) learned from a large-scale text corpus, and they adopt collaborative filtering (Sarwar et al., 2001) and matrix factorization (Koren et al., 2009) for sememe prediction, which are concluded as Sememe Prediction with Word Embeddings (SPWE) and Sememe Prediction with Sememe Embeddings (SPSE) respectively. However, those methods ignore the internal information within words (e.g., the characters in Chinese words), which is also significant for word understanding, especially for words which are of low-frequency or do not appear in the corpus at all. In this paper, we take Chinese as an example and explore methods of taking full advantage of both external and internal information of words for sememe prediction.

In Chinese, words are composed of one or multiple characters, and most characters have corresponding semantic meanings. As shown by Yin (1984), more than 90% of Chinese characters in modern Chinese corpora are morphemes. Chinese words can be divided into single-morpheme words and compound words, where compound words account for a dominant proportion. The meanings of compound words are closely related to their internal characters as shown in Fig. 1. Taking a compound word “铁匠” (ironsmith) for instance, it consists of two Chinese characters: “铁” (iron) and “匠” (craftsman), and the semantic meaning of “铁匠” can be inferred from the combination of its two characters (*iron + craftsman* → *ironsmith*). Even for some single-morpheme words, their semantic meanings may also be deduced from their characters. For example, both characters of the single-morpheme word “徘徊” (hover) represent the meaning of “hover” or “linger”. Therefore, it is intuitive to take the internal character information into consideration for sememe prediction.

In this paper, we propose a novel framework for Character-enhanced Sememe Prediction (CSP), which leverages both internal character information and external context for sememe prediction.

CSP predicts the sememe candidates for a target word from its word embedding and the corresponding character embeddings. Specifically, we follow SPWE and SPSE as introduced by Xie et al. (2017) to model external information and propose Sememe Prediction with Word-to-Character Filtering (SPWCF) and Sememe Prediction with Character and Sememe Embeddings (SPCSE) to model internal character information. In our experiments, we evaluate our models on the task of sememe prediction using HowNet. The results show that CSP achieves state-of-the-art performance and stays robust for low-frequency words.

To summarize, the key contributions of this work are as follows: (1) To the best of our knowledge, this work is the first to consider the internal information of characters for sememe prediction. (2) We propose a sememe prediction framework considering both external and internal information, and show the effectiveness and robustness of our models on a real-world dataset.

2 Related Work

Knowledge Bases. Knowledge Bases (KBs), aiming to organize human knowledge in structural forms, are playing an increasingly important role as infrastructural facilities of artificial intelligence and natural language processing. KBs rely on manual efforts (Bollacker et al., 2008), automatic extraction (Auer et al., 2007), manual evaluation (Suchanek et al., 2007), automatic completion and alignment (Bordes et al., 2013; Toutanova et al., 2015; Zhu et al., 2017) to build, verify and enrich their contents. WordNet (Miller, 1995) and BabelNet (Navigli and Ponzetto, 2012) are the representative of linguist KBs, where words of similar meanings are grouped to form thesaurus (Nastase and Szpakowicz, 2001). Apart from other linguistic KBs, sememe KBs such as HowNet (Dong and Dong, 2006) can play a significant role in understanding the semantic meanings of concepts in human languages and are favorable for various NLP tasks: information structure annotation (Gan and Wong, 2000), word sense disambiguation (Gan et al., 2002), word representation learning (Niu et al., 2017; Faruqui et al., 2015), and sentiment analysis (Fu et al., 2013) inter alia. Hence, lexical sememe prediction is an important task to construct sememe KBs.

Automatic Sememe Prediction. Automatic sememe prediction is proposed by Xie et al. (2017).

For this task, they propose SPWE and SPSE, which are inspired by collaborative filtering (Sarwar et al., 2001) and matrix factorization (Koren et al., 2009) respectively. SPWE recommends the sememes of those words that are close to the unlabelled word in the embedding space. SPSE learns sememe embeddings by matrix factorization (Koren et al., 2009) within the same embedding space of words, and it then recommends the most relevant sememes to the unlabelled word in the embedding space. In these methods, word embeddings are learned based on external context information (Pennington et al., 2014; Mikolov et al., 2013) on large-scale text corpus. These methods do not exploit internal information of words, and fail to handle low-frequency words and out-of-vocabulary words. In this paper, we propose to incorporate internal information for lexical sememe prediction.

Subword and Character Level NLP. Subword and character level NLP models the internal information of words, which is especially useful to address the out-of-vocabulary (OOV) problem. Morphology is a typical research area of subword level NLP. Subword level NLP has also been widely considered in many NLP applications, such as keyword spotting (Narasimhan et al., 2014), parsing (Seeker and Çetinoğlu, 2015), machine translation (Dyer et al., 2010), speech recognition (Creutz et al., 2007), and paradigm completion (Sutskever et al., 2014; Bahdanau et al., 2015; Cotterell et al., 2016a; Kann et al., 2017; Jin and Kann, 2017). Incorporating subword information for word embeddings (Bojanowski et al., 2017; Cotterell et al., 2016b; Chen et al., 2015; Wieting et al., 2016; Yin et al., 2016) facilitates modeling rare words and can improve the performance of several NLP tasks to which the embeddings are applied. Besides, people also consider character embeddings which have been utilized in Chinese word segmentation (Sun et al., 2014).

The success of previous work verifies the feasibility of utilizing internal character information of words. We design our framework for lexical sememe prediction inspired by these methods.

3 Background and Notation

In this section, we first introduce the organization of *sememes*, *senses* and *words* in HowNet. Then we offer a formal definition of lexical sememe prediction and develop our notation.

3.1 Sememes, Senses and Words in HowNet

HowNet provides sememe annotations for Chinese words, where each word is represented as a hierarchical tree-like sememe structure. Specifically, a word in HowNet may have various *senses*, which respectively represent the semantic meanings of the word in the real world. Each *sense* is defined as a hierarchical structure of *sememes*. For instance, as shown in the right part of Fig. 1, the word “铁匠” (ironsmith) has one sense, namely *ironsmith*. The sense *ironsmith* is defined by the sememe “人” (human) which is modified by sememe “职位” (occupation), “金属” (metal) and “工” (industrial). In HowNet, linguists use about 2,000 sememes to describe more than 100,000 words and phrases in Chinese with various combinations and hierarchical structures.

3.2 Formalization of the Task

In this paper, we focus on the relationships between the *words* and the *sememes*. Following the settings of Xie et al. (2017), we simply ignore the senses and the hierarchical structure of sememes, and we regard the sememes of all senses of a word together as the sememe set of the word.

We now introduce the notation used in this paper. Let $G = (W, S, T)$ denotes the sememe KB, where $W = \{w_1, w_2, \dots, w_{|W|}\}$ is the set of words, S is the set of sememes, and $T \subseteq W \times S$ is the set of relation pairs between words and sememes. We denote the Chinese character set as C , with each word $w_i \in C^+$. Each word w has its sememe set $S_w = \{s | (w, s) \in T\}$. Take the word “铁匠” (ironsmith) for example, the sememe set $S_{\text{铁匠 (ironsmith)}}$ consists of “人” (human), “职位” (occupation), “金属” (metal) and “工” (industrial).

Given a word $w \in C^+$, the task of lexical sememe prediction aims to predict the corresponding $P(s|w)$ of sememes in S to recommend them to w .

4 Methodology

In this section, we present our framework for lexical sememe prediction (SP). For each unlabelled word, our framework aims to recommend the most appropriate sememes based on the internal and external information. Because of introducing character information, our framework can work for both high-frequency and low-frequency words.

Our framework is the ensemble of two parts: sememe prediction with internal information (i.e., *internal* models), and sememe prediction with external information (i.e., *external* models). Explicitly, we adopt SPWE, SPSE, and their ensemble (Xie et al., 2017) as *external* models, and we take SPWCF, SPCSE, and their ensemble as *internal* models.

In the following sections, we first introduce SPWE and SPSE. Then, we show the details of SPWCF and SPCSE. Finally, we present the method of model ensembling.

4.1 SP with External Information

SPWE and SPSE are introduced by Xie et al. (2017) as the state of the art for sememe prediction. These methods represent word meanings with embeddings learned from external information, and apply the ideas of collaborative filtering and matrix factorization in recommendation systems for sememe predication.

SP with Word Embeddings (SPWE) is based on the assumption that similar words should have similar sememes. In SPWE, the similarity of words are measured by cosine similarity. The score function $P(s_j|w)$ of sememe s_j given a word w is defined as:

$$P(s_j|w) \sim \sum_{w_i \in W} \cos(\mathbf{w}, \mathbf{w}_i) \cdot \mathbf{M}_{ij} \cdot c^{r_i}, \quad (1)$$

where \mathbf{w} and \mathbf{w}_i are pre-trained word embeddings of words w and w_i . $\mathbf{M}_{ij} \in \{0, 1\}$ indicates the annotation of sememe s_j on word w_i , where $\mathbf{M}_{ij} = 1$ indicates the word $s_j \in S_{w_i}$ and otherwise is not. r_i is the descend cosine word similarity rank between w and w_i , and $c \in (0, 1)$ is a hyper-parameter.

SP with Sememe Embeddings (SPSE) aims to map sememes into the same low-dimensional space of the word embeddings to predict the semantic correlations of the sememes and the words. This method learns two embeddings \mathbf{s} and $\bar{\mathbf{s}}$ for each sememe by solving matrix factorization with the loss function defined as:

$$\mathcal{L} = \sum_{w_i \in W, s_j \in S} (\mathbf{w}_i \cdot (\mathbf{s}_j + \bar{\mathbf{s}}_j) + \mathbf{b}_i + \mathbf{b}'_j - \mathbf{M}_{ij})^2 + \lambda \sum_{s_j, s_k \in S} (\mathbf{s}_j \cdot \bar{\mathbf{s}}_k - \mathbf{C}_{jk})^2, \quad (2)$$

where \mathbf{M} is the same matrix used in SPWE. \mathbf{C} indicates the correlations between sememes, in

which \mathbf{C}_{jk} is defined as the point-wise mutual information $\text{PMI}(s_j, s_k)$. The sememe embeddings are learned by factorizing the word-sememe matrix \mathbf{M} and the sememe-sememe matrix \mathbf{C} synchronously with fixed word embeddings. \mathbf{b}_i and \mathbf{b}'_j denote the bias of w_i and s_j , and λ is a hyper-parameter. Finally, the score of sememe s_j given a word w is defined as:

$$P(s_j|w) \sim \mathbf{w} \cdot (\mathbf{s}_j + \bar{\mathbf{s}}_j). \quad (3)$$

4.2 SP with Internal Information

We design two methods for sememe prediction with only internal character information without considering contexts as well as pre-trained word embeddings.

4.2.1 SP with Word-to-Character Filtering (SPWCF)

Inspired by collaborative filtering (Sarwar et al., 2001), we propose to recommend sememes for an unlabelled word according to its similar words based on internal information. Instead of using pre-trained word embeddings, we consider words as *similar* if they contain the same characters at the same positions.

In Chinese, the meaning of a character may vary according to its position within a word (Chen et al., 2015). We consider three positions within a word: *Begin*, *Middle*, and *End*. For example, as shown in Fig. 2, the character at the *Begin* position of the word “火车站” (railway station) is “火” (fire), while “车” (vehicle) and “站” (station) are at the *Middle* and *End* position respectively. The character “站” usually means *station* when it is at the *End* position, while it usually means *stand* at the *Begin* position like in “站立” (stand), “站岗哨兵” (standing guard) and “站起来” (stand up).



Figure 2: An example of the position of characters in a word.

Formally, for a word $w = c_1 c_2 \dots c_{|w|}$, we define $\pi_B(w) = \{c_1\}$, $\pi_M(w) = \{c_2, \dots, c_{|w|-1}\}$, $\pi_E(w) = \{c_{|w|}\}$, and

$$P_p(s_j|c) \sim \frac{\sum_{w_i \in W \wedge c \in \pi_p(w_i)} \mathbf{M}_{ij}}{\sum_{w_i \in W \wedge c \in \pi_p(w_i)} |S_{w_i}|}, \quad (4)$$

that represents the score of a sememe s_j given a character c and a position p , where π_p may be π_B , π_M , or π_E . \mathbf{M} is the same matrix used in Eq. (1). Finally, we define the score function $P(s_j|w)$ of sememe s_j given a word w as:

$$P(s_j|w) \sim \sum_{p \in \{B, M, E\}} \sum_{c \in \pi_p(w)} P_p(s_j|c). \quad (5)$$

SPWCF is a simple and efficient method. It performs well because compositional semantics are pervasive in Chinese compound words, which makes it straightforward and effective to find similar words according to common characters.

4.2.2 SP with Character and Sememe Embeddings (SPCSE)

The method Sememe Prediction with Word-to-Character Filtering (SPWCF) can effectively recommend the sememes that have strong correlations with characters. However, just like SPWE, it ignores the relations between sememes. Hence, inspired by SPSE, we propose Sememe Prediction with Character and Sememe Embeddings (SPCSE) to take the relations between sememes into account. In SPCSE, we instead learn the sememe embeddings based on internal character information, then compute the semantic distance between sememes and words for prediction.

Inspired by GloVe (Pennington et al., 2014) and SPSE, we adopt matrix factorization in SPCSE, by decomposing the word-sememe matrix and the sememe-sememe matrix simultaneously. Instead of using pre-trained word embeddings in SPSE, we use pre-trained character embeddings in SPCSE. Since the ambiguity of characters is stronger than that of words, multiple embeddings are learned for each character (Chen et al., 2015). We select the most representative character and its embedding to represent the word meaning. Because low-frequency characters are much rare than those low-frequency words, and even low-frequency words are usually composed of common characters, it is feasible to use pre-trained character embeddings to represent rare words. During factorizing the word-sememe matrix, the character embeddings are fixed.

We set N_e as the number of embeddings for each character, and each character c has N_e embeddings $\mathbf{c}^1, \dots, \mathbf{c}^{N_e}$. Given a word w and a sememe s , we select the embedding of a character of w closest to the sememe embedding by cosine distance as the representation of the word w ,

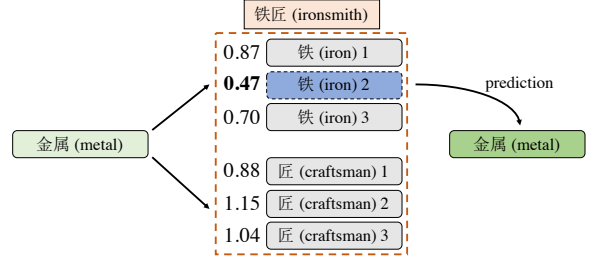


Figure 3: An example of adopting multiple-prototype character embeddings. The numbers are the cosine distances. The sememe “金属” (metal) is the closest to one embedding of “铁” (iron).

as shown in Fig. 3. Specifically, given a word $w = c_1 \dots c_{|w|}$ and a sememe s_j , we define

$$\hat{k}, \hat{r} = \arg \min_{k, r} [1 - \cos(\mathbf{c}_k^r, (\mathbf{s}'_j + \bar{\mathbf{s}}'_j))], \quad (6)$$

where \hat{k} and \hat{r} indicate the indices of the character and its embedding closest to the sememe s_j in the semantic space. With the same word-sememe matrix \mathbf{M} and sememe-sememe correlation matrix \mathbf{C} in Eq. (2), we learn the sememe embeddings with the loss function:

$$\mathcal{L} = \sum_{w_i \in W, s_j \in S} \left(\mathbf{c}_{\hat{k}}^{\hat{r}} \cdot (\mathbf{s}'_j + \bar{\mathbf{s}}'_j) + \mathbf{b}_k^c + \mathbf{b}_j'' - \mathbf{M}_{ij} \right)^2 + \lambda' \sum_{s_j, s_q \in S} (\mathbf{s}'_j \cdot \bar{\mathbf{s}}'_q - \mathbf{C}_{jq})^2, \quad (7)$$

where \mathbf{s}'_j and $\bar{\mathbf{s}}'_j$ are the sememe embeddings for sememe s_j , and $\mathbf{c}_{\hat{k}}^{\hat{r}}$ is the embedding of the character that is the closest to sememe s_j within w_i . Note that, as the characters and the words are not embedded into the same semantic space, we learn new sememe embeddings instead of using those learned in SPSE, hence we use different notations for the sake of distinction. \mathbf{b}_k^c and \mathbf{b}_j'' denote the biases of c_k and s_j , and λ' is the hyper-parameter adjusting the two parts. Finally, the score function of word $w = c_1 \dots c_{|w|}$ is defined as:

$$P(s_j|w) \sim \mathbf{c}_{\hat{k}}^{\hat{r}} \cdot (\mathbf{s}'_j + \bar{\mathbf{s}}'_j). \quad (8)$$

4.3 Model Ensembling

SPWCF / SPCSE and SPWE / SPSE take different sources of information as input, which means that they have different characteristics: SPWCF / SPCSE only have access to internal information, while SPWE / SPSE can only make use of external

information. On the other hand, just like the difference between SPWE and SPSE, SPWCF originates from collaborative filtering, whereas SPCSE uses matrix factorization. All of those methods have in common that they tend to recommend the sememes of *similar* words, but they diverge in their interpretation of *similar*.

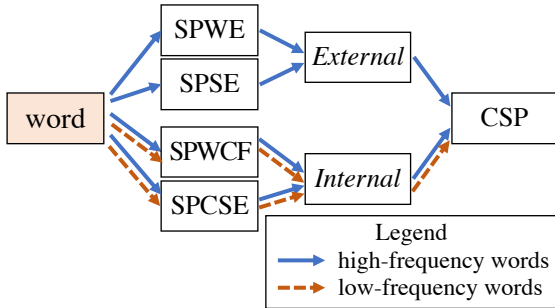


Figure 4: The illustration of model ensembling.

Hence, to obtain better prediction performance, it is necessary to combine these models. We denote the ensemble of SPWCF and SPCSE as the *internal* model, and we denote the ensemble of SPWE and SPSE as the *external* model. The ensemble of the *internal* and the *external* models is our novel framework CSP. In practice, for words with reliable word embeddings, i.e., high-frequency words, we can use the integration of the *internal* and the *external* models; for words with extremely low frequencies (e.g., having no reliable word embeddings), we can just use the *internal* model and ignore the *external* model, because the external information is noise in this case. Fig. 4 shows model ensembling in different scenarios. For the sake of comparison, we use the integration of SPWCF, SPCSE, SPWE, and SPSE as CSP in our all experiments. In this paper, two models are integrated by simple weighted addition.

5 Experiments

In this section, we evaluate our models on the task of sememe prediction. Additionally, we analyze the performance of different methods for various word frequencies. We also execute an elaborate case study to demonstrate the mechanism of our methods and the advantages of using internal information.

5.1 Dataset

We use the human-annotated sememe KB HowNet for sememe prediction. In HowNet, 103,843

words are annotated with 212,539 senses, and each sense is defined as a hierarchical structure of sememes. There are about 2,000 sememes in HowNet. However, the frequencies of some sememes in HowNet are very low, such that we consider them unimportant and remove them. Our final dataset contains 1,400 sememes. For learning the word and character embeddings, we use the Sogou-T corpusⁱⁱ (Liu et al., 2012), which contains 2.7 billion words.

5.2 Experimental Settings

In our experiments, we evaluate SPWCF, SPCSE, and SPWCF + SPCSE which only use internal information, and the ensemble framework CSP which uses both internal and external information for sememe prediction. We use the state-of-the-art models from Xie et al. (2017) as our baselines. Additionally, we use the SPWE model with word embeddings learned by fastText (Bojanowski et al., 2017) that considers both internal and external information as a baseline.

For the convenience of comparison, we select 60,000 high-frequency words in Sogou-T corpus from HowNet. We divide the 60,000 words into train, dev, and test sets of size 48,000, 6,000, and 6,000, respectively, and we keep them fixed throughout all experiments except for Section 5.4. In Section 5.4, we utilize the same train and dev sets, but use other words from HowNet as the test set to analyze the performance of our methods for different word frequency scenarios. We select the hyper-parameters on the dev set for all models including the baselines and report the evaluation results on the test set.

We set the dimensions of the word, sememe, and character embeddings to be 200. The word embeddings are learned by GloVe (Pennington et al., 2014). For the baselines, in SPWE, the hyper-parameter c is set to 0.8, and the model considers no more than $K = 100$ nearest words. We set the probability of decomposing zero elements in the word-sememe matrix in SPSE to be 0.5%. λ in Eq. (2) is 0.5. The model is trained for 20 epochs, and the initial learning rate is 0.01, which decreases through iterations. For fastText, we use skip-gram with hierarchical softmax to learn word embeddings, and we set the minimum length of character n-grams to be 1 and the maximum length

ⁱⁱ Sogou-T corpus is provided by Sogou Inc., a Chinese commercial search engine company. <https://www.sogou.com/labs/resource/t.php>

of character n-grams to be 2. For model ensembling, we use $\frac{\lambda_{SPWE}}{\lambda_{SPSE}} = 2.1$ as the addition weight.

For SPCSE, we use Cluster-based Character Embeddings (Chen et al., 2015) to learn pre-trained character embeddings, and we set N_e to be 3. We set λ' in Eq. (7) to be 0.1, and the model is trained for 20 epochs. The initial learning rate is 0.01 and decreases during training as well. Since generally each character can relate to about 15 - 20 sememes, we set the probability of decomposing zero elements in the word-sememe matrix in SPCSE to be 2.5%. The ensemble weight of SPWCF and SPCSE $\frac{\lambda_{SPWCF}}{\lambda_{SPCSE}} = 4.0$. For better performance of the final ensemble model CSP, we set $\lambda = 0.1$ and $\frac{\lambda_{SPWE}}{\lambda_{SPSE}} = 0.3125$, though 0.5 and 2.1 are the best for SPSE and SPWE + SPSE. Finally, we choose $\frac{\lambda_{internal}}{\lambda_{external}} = 1.0$ to integrate the *internal* and *external* models.

5.3 Sememe Prediction

5.3.1 Evaluation Protocol

The task of sememe prediction aims to recommend appropriate sememes for unlabelled words. We cast this as a multi-label classification task, and adopt mean average precision (MAP) as the evaluation metric. For each unlabelled word in the test set, we rank all sememe candidates with the scores given by our models as well as baselines, and we report the MAP results. The results are reported on the test set, and the hyper-parameters are tuned on the dev set.

5.3.2 Experiment Results

The evaluation results are shown in Table 1. We can observe that:

Method	MAP
SPSE	0.411
SPWE	0.565
SPWE+SPSE	0.577
SPWCF	0.467
SPCSE	0.331
SPWCF + SPCSE	0.483
SPWE + fastText	0.531
CSP	0.654

Table 1: Evaluation results on sememe prediction. The result of SPWCF + SPCSE is bold for comparing with other methods (SPWCF and SPCSE) which use only internal information.

(1) Considerable improvements are obtained via model ensembling, and the CSP model achieves state-of-the-art performance. CSP combines the internal character information with the external context information, which significantly and consistently improves performance on sememe prediction. Our results confirm the effectiveness of a combination of internal and external information for sememe prediction; since different models focus on different features of the inputs, the ensemble model can absorb the advantages of both methods.

(2) The performance of SPWCF + SPCSE is better than that of SPSE, which means using only internal information could already give good results for sememe prediction as well. Moreover, in *internal* models, SPWCF performs much better than SPCSE, which also implies the strong power of collaborative filtering.

(3) The performance of SPWCF + SPCSE is worse than SPWE + SPSE. This indicates that it is still difficult to figure out the semantic meanings of a word without contextual information, due to the ambiguity and meaning vagueness of internal characters. Moreover, some words are not compound words (e.g., single-morpheme words or transliterated words), whose meanings can hardly be inferred directly by their characters. In Chinese, internal character information is just partial knowledge. We present the results of SPWCF and SPCSE merely to show the capability to use the internal information in isolation. In our case study, we will demonstrate that *internal* models are powerful for low-frequency words, and can be used to predict senses that do not appear in the corpus.

5.4 Analysis on Different Word Frequencies

To verify the effectiveness of our models on different word frequencies, we incorporate the remaining words in HowNetⁱⁱⁱ into the test set. Since the remaining words are low-frequency, we mainly focus on words with long-tail distribution. We count the number of occurrences in the corpus for each word in the test set and group them into eight categories by their frequency. The evaluation results are shown in Table 2, from which we can observe that:

ⁱⁱⁱ In detail, we do not use the numeral words, punctuations, single-character words, the words do not appear in Sogou-T corpus (because they need to appear at least for one time to get the word embeddings), and foreign abbreviations.

word frequency occurrences	≤ 50	51– 100	101 – 1,000	1,001 – 5,000	5,001 – 10,000	10,001 – 30,000	>30,000
	8537	4868	3236	2036	663	753	686
SPWE	0.312	0.437	0.481	0.558	0.549	0.556	0.509
SPSE	0.187	0.273	0.339	0.409	0.407	0.424	0.386
SPWE + SPSE	0.284	0.414	0.478	0.556	0.548	0.554	0.511
SPWCF	0.456	0.414	0.400	0.443	0.462	0.463	0.479
SPCSE	0.309	0.291	0.286	0.312	0.339	0.353	0.342
SPWCF + SPCSE	0.467	0.437	0.418	0.456	0.477	0.477	0.494
SPWE + fastText	0.495	0.472	0.462	0.520	0.508	0.499	0.490
CSP	0.527	0.555	0.555	0.626	0.632	0.641	0.624

Table 2: MAP scores on sememe prediction with different word frequencies.

words	models	Top 5 sememes
钟表匠 (clockmaker)	internal	人(human), 职位(occupation), 部件(part), 时间(time), 告诉(tell)
	external	人(human), 专(ProperName), 地方(place), 欧洲(Europe), 政(politics)
	ensemble	人(human), 职位(occupation), 告诉(tell), 时间(time), 用具(tool)
奥斯卡 (Oscar)	internal	专(ProperName), 地方(place), 市(city), 人(human), 国都(capital)
	external	奖励(reward), 艺(entertainment), 专(ProperName), 用具(tool), 事情(fact)
	ensemble	专(ProperName), 奖励(reward), 艺(entertainment), 著名(famous), 地方(place)

Table 3: Examples of sememe prediction. For each word, we present the top 5 sememes predicted by the *internal* model, *external* model and the final ensemble model (CSP). Bold sememes are correct.

(1) The performances of SPSE, SPWE, and SPWE + SPSE decrease dramatically with low-frequency words compared to those with high-frequency words. On the contrary, the performances of SPWCF, SPCSE, and SPWCF + SPCSE, though weaker than that on high-frequency words, is not strongly influenced in the long-tail scenario. The performance of CSP also drops since CSP also uses external information, which is not sufficient with low-frequency words. These results show that the word frequencies and the quality of word embeddings can influence the performance of sememe prediction methods, especially for *external* models which mainly concentrate on the word itself. However, the *internal* models are more robust when encountering long-tail distributions. Although words do not need to appear too many times for learning good word embeddings, it is still hard for *external models* to recommend sememes for low-frequency words. While since *internal* models do not use external word embeddings, they can still work in such scenario. As for the performance on high-frequency words, since these words are used widely, the ambiguity of high-frequency words is thus much stronger, while the *internal* models are still stable for high-frequency words.

(2) The results also indicate that even low-frequency words in Chinese are mostly composed of common characters, and thus it is possible

to utilize internal character information for sememe prediction on words with long-tail distribution (even on those new words that never appear in the corpus). Moreover, the stability of the MAP scores given by our methods on various word frequencies also reflects the reliability and universality of our models in real-world sememe annotations in HowNet. We will give detailed analysis in our case study.

5.5 Case Study

The results of our main experiments already show the effectiveness of our models. In this case study, we further investigate the outputs of our models to confirm that character-level knowledge is truly incorporated into sememe prediction.

In Table 3, we demonstrate the top 5 sememes for “钟表匠” (clockmaker) and “奥斯卡” (Oscar, i.e., the Academy Awards). “钟表匠” (clockmaker) is a typical compound word, while “奥斯卡” (Oscar) is a transliterated word. For each word, the top 5 results generated by the internal model (SPWCF + SPCSE), the *external* model (SPWE + SPSE) and the ensemble model (CSP) are listed.

The word “钟表匠” (clockmaker) is composed of three characters: “钟” (bell, clock), “表” (clock, watch) and “匠” (craftsman). Humans can intuitively conclude that *clock + craftsman* → *clockmaker*. However, the *external* model does not per-

form well for this example. If we investigate the word embedding of “钟表匠” (clockmaker), we can know why this method recommends these unreasonable sememes. The closest 5 words in the train set to “钟表匠” (clockmaker) by cosine similarity of their embeddings are: “瑞士” (Switzerland), “卢梭” (Jean-Jacques Rousseau), “鞋匠” (cobbler), “发明家” (inventor) and “奥地利人” (Austrian). Note that none of these words are directly relevant to *bells*, *clocks* or *watches*. Hence, the sememes “时间” (time), “告诉” (tell), and “用具” (tool) cannot be inferred by those words, even though the correlations between sememes are introduced by SPSE. In fact, those words are related to *clocks* in an indirect way: Switzerland is famous for watch industry; Rousseau was born into a family that had a tradition of watchmaking; cobbler and inventor are two kinds of occupations as well. With the above reasons, those words usually co-occur with “钟表匠” (clockmaker), or usually appear in similar contexts as “钟表匠” (clockmaker). It indicates that related word embeddings as used in an *external* model do not always recommend related sememes.

The word “奥斯卡” (Oscar) is created by the pronunciation of *Oscar*. Therefore, the meaning of each character in “奥斯卡” (Oscar) is unrelated to the meaning of the word. Moreover, the characters “奥”, “斯”, and “卡” are common among transliterated words, thus the *internal* method recommends “专” (ProperName) and “地方” (place), etc., since many transliterated words are proper nouns or place names.

6 Conclusion and Future Work

In this paper, we introduced character-level internal information for lexical sememe prediction in Chinese, in order to alleviate the problems caused by the exclusive use of external information. We proposed a Character-enhanced Sememe Prediction (CSP) framework which integrates both internal and external information for lexical sememe prediction and proposed two methods for utilizing internal information. We evaluated our CSP framework on the classical manually annotated sememe KB HowNet. In our experiments, our methods achieved promising results and outperformed the state of the art on sememe prediction, especially for low-frequency words.

We will explore the following research directions in the future: (1) Concepts in HowNet are an-

notated with hierarchical structures of senses and sememes, but those are not considered in this paper. In the future, we will take structured annotations into account. (2) It would be meaningful to take more information into account for blending external and internal information and design more sophisticated methods. (3) Besides Chinese, many other languages have rich subword-level information. In the future, we will explore methods of exploiting internal information in other languages. (4) We believe that sememes are universal for all human languages. We will explore a general framework to recommend and utilize sememes for other NLP tasks.

Acknowledgments

This research is part of the NEX++ project, supported by the National Research Foundation, Prime Minister’s Office, Singapore under its IRC@Singapore Funding Initiative. This work is also supported by the National Natural Science Foundation of China (NSFC No. 61661146007 and 61572273) and the research fund of Tsinghua University-Tencent Joint Laboratory for Internet Innovation Technology. Hao Zhu is supported by Tsinghua University Initiative Scientific Research Program. We would like to thank Katharina Kann, Shen Jin, and the anonymous reviewers for their helpful comments.

References

- Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. Dbpedia: A nucleus for a web of open data. In *Proceedings of ISWC*, pages 722–735.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR*.
- Leonard Bloomfield. 1926. A set of postulates for the science of language. *Language*, 2(3):153–164.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *TACL*, 5:135–146.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: A collaboratively created graph database for structuring human knowledge. In *Proceedings of SIGMOD*, pages 1247–1250.

- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Proceedings of NIPS*, pages 2787–2795.
- Xinxiong Chen, Lei Xu, Zhiyuan Liu, Maosong Sun, and Huan-Bo Luan. 2015. Joint learning of character and word embeddings. In *Proceedings of IJCAI*, pages 1236–1242.
- Ryan Cotterell, Christo Kirov, John Sylak-Glassman, David Yarowsky, Jason Eisner, and Mans Hulden. 2016a. The SIGMORPHON 2016 shared task—morphological reinflection. In *Proceedings of SIGMORPHON*, pages 10–22.
- Ryan Cotterell, Hinrich Schütze, and Jason Eisner. 2016b. Morphological smoothing and extrapolation of word embeddings. In *Proceedings of ACL*, pages 1651–1660.
- Mathias Creutz, Teemu Hirsimäki, Mikko Kurimo, Antti Puurula, Janne Pytkönen, Vesa Siivola, Matti Varjokallio, Ebru Arisoy, Murat Saraçlar, and Andreas Stolcke. 2007. Analysis of morph-based speech recognition and the modeling of out-of-vocabulary words across languages. In *Proceedings of HLT-NAACL*, pages 380–387.
- Zhendong Dong and Qiang Dong. 2006. *HowNet and the computation of meaning*. World Scientific.
- Xiangyu Duan, Jun Zhao, and Bo Xu. 2007. Word sense disambiguation through sememe labeling. In *Proceedings of IJCAI*, pages 1594–1599.
- Chris Dyer, Jonathan Weese, Hendra Setiawan, Adam Lopez, Ferhan Ture, Vladimir Eidelman, Juri Ganitkevitch, Phil Blunsom, and Philip Resnik. 2010. cdec: A decoder, alignment, and learning framework for finite-state and context-free translation models. In *Proceedings of the ACL 2010 System Demonstrations*, pages 7–12.
- Manaal Faruqui, Jesse Dodge, Sujay Kumar Jauhar, Chris Dyer, Eduard Hovy, and Noah A. Smith. 2015. Retrofitting word vectors to semantic lexicons. In *Proceedings of HLT-NAACL*, pages 1606–1615.
- Xianghua Fu, Liu Guo, Guo Yanyan, and Wang Zhiqiang. 2013. Multi-aspect sentiment analysis for Chinese online social reviews based on topic modeling and HowNet lexicon. *Knowledge-Based Systems*, 37:186–195.
- Kok-Wee Gan, Chi-Yung Wang, and Brian Mak. 2002. Knowledge-based sense pruning using the HowNet: an alternative to word sense disambiguation. In *Proceedings of ISCSLP*.
- Kok Wee Gan and Ping Wai Wong. 2000. Annotating information structures in Chinese texts using HowNet. In *Proceedings of The Second Chinese Language Processing Workshop*, pages 85–92.
- Minlie Huang, Borui Ye, Yichen Wang, Haiqiang Chen, Junjun Cheng, and Xiaoyan Zhu. 2014. New word detection for sentiment analysis. In *Proceedings of ACL*, pages 531–541.
- Huiming Jin and Katharina Kann. 2017. Exploring cross-lingual transfer of morphological knowledge in sequence-to-sequence models. In *Proceedings of SCLeM*, pages 70–75.
- Katharina Kann, Ryan Cotterell, and Hinrich Schütze. 2017. One-shot neural cross-lingual transfer for paradigm completion. In *Proceedings of ACL*, pages 1993–2003.
- Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer*, 42(8).
- Yiqun Liu, Fei Chen, Weize Kong, Huijia Yu, Min Zhang, Shaoping Ma, and Liyun Ru. 2012. Identifying web spam with the wisdom of the crowds. *ACM Transactions on the Web*, 6(1):2:1–2:30.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of NIPS*, pages 3111–3119.
- George A Miller. 1995. WordNet: A lexical database for English. *Communications of the ACM*, 38(11):39–41.
- Karthik Narasimhan, Damianos Karakos, Richard Schwartz, Stavros Tsakalidis, and Regina Barzilay. 2014. Morphological segmentation for keyword spotting. In *Proceedings of EMNLP*, pages 880–885.
- Vivi Nastase and Stan Szpakowicz. 2001. Word sense disambiguation in Roget’s thesaurus using WordNet. In *Proceedings of the Workshop on WordNet and Other Lexical Resources: Applications, Extensions and Customizations*.
- Roberto Navigli and Simone Paolo Ponzetto. 2012. BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, 193:217–250.
- Yilin Niu, Ruobing Xie, Zhiyuan Liu, and Maosong Sun. 2017. Improved word representation learning with sememes. In *Proceedings of ACL*, pages 2049–2058.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of EMNLP*, pages 1532–1543.
- Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. 2001. Item-based collaborative filtering recommendation algorithms. In *Proceedings of WWW*, pages 285–295.

- Wolfgang Seeker and Özlem Çetinoğlu. 2015. A graph-based lattice dependency parser for joint morphological segmentation and syntactic analysis. *TACL*, 3:359–373.
- Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. 2007. Yago: A core of semantic knowledge. In *Proceedings of WWW*, pages 697–706.
- Yaming Sun, Lei Lin, Nan Yang, Zhenzhou Ji, and Xiaolong Wang. 2014. Radical-enhanced Chinese character embedding. In *Proceedings of ICONIP*, pages 279–286.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of NIPS*, pages 3104–3112.
- Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoi-fung Poon, Pallavi Choudhury, and Michael Gamon. 2015. Representing text for joint embedding of text and knowledge bases. In *Proceedings of EMNLP*, pages 1499–1509.
- John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2016. Charagram: Embedding words and sentences via character n-grams. In *Proceedings of EMNLP*, pages 1504–1515.
- Yunqing Xia, Taotao Zhao, Jianmin Yao, and Peng Jin. 2011. Measuring Chinese-English cross-lingual word similarity with HowNet and parallel corpus. In *Proceedings of CICLing*, pages 221–233. Springer.
- Ruobing Xie, Xingchi Yuan, Zhiyuan Liu, and Maosong Sun. 2017. Lexical sememe prediction via word embeddings and matrix factorization. In *Proceedings of IJCAI*, pages 4200–4206.
- Binyong Yin. 1984. Quantitative research on Chinese morphemes. *Studies of the Chinese Language*, 5:338–347.
- Rongchao Yin, Quan Wang, Peng Li, Rui Li, and Bin Wang. 2016. Multi-granularity Chinese word embedding. In *Proceedings of EMNLP*, pages 981–986.
- Hao Zhu, Ruobing Xie, Zhiyuan Liu, and Maosong Sun. 2017. Iterative entity alignment via joint knowledge embeddings. In *Proceedings of IJCAI*, pages 4258–4264.