

# Effective and Efficient Word-type aware Chinese Lexical Simplification

Anonymous ACL submission

## Abstract

In this paper, we address the task of Chinese lexical simplification (CLS), which aims to replace complex words in a given sentence with simpler alternatives of equivalent meaning. We propose an effective and efficient CLS method that combines small and large models in a complementary way based on the type of complex words. Specifically, we analyze the strengths and weaknesses of small models and ChatGPT. We find that ChatGPT performs well in simplifying in-dictionary common words and Chinese idioms, while small models struggle with them. Therefore, we propose an automatic knowledge distillation approach to fine-tune small models with in-dictionary words-oriented training data generated by ChatGPT. On the other hand, we find that both small models and ChatGPT have difficulties with out-of-dictionary (OOD) words. To address this issue, we use a retrieval-based interpretation augmentation strategy to enrich the input with relevant information obtained from external sources. With this strategy, both small models and ChatGPT can significantly improve their performance in simplifying OOD words. Finally, we introduce a simple controller that selects the best model or tool for each complex word according to its type. This hybrid approach can balance performance and cost and achieve better results than any single model.

## 1 Introduction

Lexical Simplification (LS) is the task of replacing complex words in a sentence with simpler alternatives while preserving their structure and original meaning. This task can improve the readability of the text to benefit a wide range of people, such as students (De Belder and Moens, 2010), non-native speakers (Paetzold and Specia, 2016), and individuals with cognitive impairments (Feng, 2009; Saggion, 2017). However, it is a challenging task that requires both linguistic knowledge and contextual awareness.

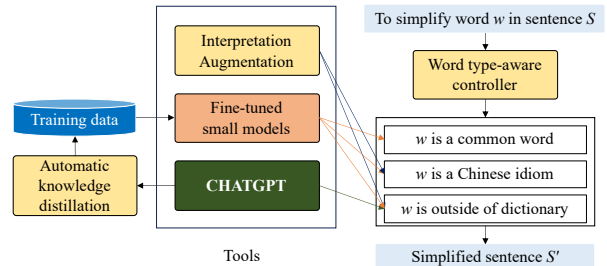


Figure 1: The general framework of the proposed word type-aware Chinese lexical simplification method.

This paper focuses on Chinese lexical simplification (CLS). One big barrier for CLS is the lack of enough training data. So recent work focuses on unsupervised methods based on pre-trained language models (PLMs), e.g., the state-of-the-art CLS system BERT-LS (Qiang et al., 2021) generates substitution candidates based on the pre-trained masked language model (MLM) BERT. Despite its simplicity, the model cannot fully understand the task, resulting in conservative word substitution and performance bottleneck.

We observe that the recent large generative pre-trained language models, such as ChatGPT (GPT-3.5 and GPT-4 (OpenAI, 2023)), can understand the task better through task instructions and a few demonstrations, while a medium model with 6B parameters still cannot get a satisfactory performance, reflecting the decisive role of model scale. However, the cost of training, maintaining, and invoking large language models is enormous. We face a trade-off between performance and cost when choosing between small or large models.

In this paper, we aim to improve a small model by learning from and collaborating with ChatGPT. We expect the final system to achieve competitive performance compared to ChatGPT while significantly reducing the inference cost. To accomplish this goal, we present the following contributions.

First, we conduct a thorough analysis of the un-

072 supervised CLS methods based on small, medium, 122  
073 and large language models, to gain a deeper under- 123  
074 standing of the advantages and disadvantages of the 124  
075 current methods. We discover that ChatGPT has 125  
076 advantages in task understanding, reducing the loss 126  
077 of details and degree of information compared with 127  
078 small models. Linguistic resources can help small 128  
079 models obtain competitive performance for com- 129  
080 mon words in the dictionary. All the models have a 130  
081 lot of room for improvement on out-of-dictionary 131  
082 (OOD) words. 132

083 Second, we propose a knowledge distillation 133  
084 framework called **PivotKD** for fine-tuning small 134  
085 models with the in-dictionary words-oriented train- 135  
086 ing data, which are generated by ChatGPT. Piv- 136  
087 otKD samples pivot words from a dictionary, 137  
088 lets ChatGPT generate sentences containing pivot 138  
089 words, and replace them with alternatives belong- 139  
090 ing to different lexical difficulty levels in an au- 140  
091 tomatic way. Evaluation shows that fine-tuning a 141  
092 model with 700m or 6B parameters can obtain su- 142  
093 perior performance compared with ChatGPT on 143  
094 simplifying common words and Chinese idioms. 144

095 Third, we propose a retrieval-based interpreta- 145  
096 tion augmentation strategy for enhancing simplifi- 146  
097 cation on OOD words. We ask a search engine for 147  
098 the interpretation of a target complex word and use 148  
099 the interpretation for in-context learning. Based on 149  
100 this simple strategy, both ChatGPT and the fine- 150  
101 tuned small models gain large improvements in 151  
102 simplifying OOD words. 152

103 Finally, as shown in Figure 1, we propose a sim- 153  
104 ple controller that selects the best model or tool 154  
105 for each complex word according to its type. The 155  
106 hybrid approach can balance performance and cost 156  
107 and get better results than any single model. 157

## 108 2 Related Work 158

109 Lexical Simplification is the process of replacing 159  
110 complex words in a given sentence with simpler 160  
111 alternatives of equivalent meaning (Paetzold and 161  
112 Specia, 2017b). This task makes contributions to 162  
113 performing text simplification focusing on lexical 163  
114 information and has wide applications in assisting 164  
115 readers with low language proficiency, cognitive 165  
116 impairments, or disabilities. Traditionally, lexical 166  
117 simplification mainly consists of a pipeline: the 167  
118 identification of complex words, the generation of 168  
119 substitution candidates, the selection of those can- 169  
120 didates based on the context, and the ranking of the 170  
121 selected substitutes according to their simplicity. 171

Complex word identification aims to identify 122  
which word is considered complex in a sentence 123  
by a given target population (Shardlow, 2013; Yi- 124  
mam et al., 2018; Dehghan et al., 2022). In this 125  
paper, we do not pay attention to this aspect and 126  
assume that the target complex words are given by 127  
the users. Readers can refer to a recent survey for 128  
more information (North et al., 2023). 129

**Knowledge-based methods** Early lexical sim- 130  
plification research relied on lexical knowledge 131  
databases to generate substitutions (Carroll et al., 132  
1998; Drndarevic and Saggion, 2012). However, 133  
the databases are expensive to construct and update, 134  
and have a limited word coverage. 135

**Word embedding-based methods** With the advent 136  
of deep learning, semantic similarity computation 137  
based on word embeddings has become a popu- 138  
lar method for substitution generation and rank- 139  
ing (Paetzold and Specia, 2017a). But this method 140  
still suffers from the word coverage problem. 141

**PLM-based methods** Subsequently, pre-trained 142  
language models (PLMs) have been suggested for 143  
this task. For example, BERT-LS (Qiang et al., 144  
2020) proposed an unsupervised method, employ- 145  
ing BERT to generate substitutions for target com- 146  
plex words according to the encoding of the sur- 147  
rounding context. PromptLS (Vásquez-Rodríguez 148  
et al., 2022) found that fine-tuning PLMs can ob- 149  
tain better performance than the unsupervised set- 150  
ting. ConLS (Sheang et al., 2022) fine-tuned an 151  
encoder-decoder model T5 for substitution gener- 152  
ation which naturally predicts simple words with 153  
multiple tokens. One challenge of fine-tuning is 154  
the scarcity of supervised training data for some 155  
languages, such as Chinese. 156

**LLM-based methods** Recently, large language 157  
models (LLMs) such as GPT-3 have been applied 158  
for lexical simplification through prompt learning- 159  
based methods. It shows that GPT-3 can understand 160  
the task and learn to predict based on task instruc- 161  
tions and a few demonstrations and obtains good 162  
performance for the English language (Aumiller 163  
and Gertz, 2022). This indicates that LLMs already 164  
embed rich linguistic knowledge and have a strong 165  
in-context learning ability. The key is to find proper 166  
ways to guide LLMs to generate the required out- 167  
put. However, training, deploying, and applying 168  
LLMs is still very expensive. 169

This paper focuses on the Chinese language. Pre- 170  
vious methods are mostly unsupervised (Qiang 171  
et al., 2021) due to the lack of enough training 172  
data, resulting in a performance bottleneck. Our 173

174	motivation is to build a system that can effectively	Kaoshi, HSK), and each complex word has 8.51 an-	222
175	combine small and large models to keep a balance	notated simple alternatives on average as reference	223
176	between performance and cost. The following work	answers.	224
177	is also relevant to our research.		
178	<b>Knowledge distillation</b> Knowledge distillation	<b>3.2.2 Evaluation Metrics</b>	225
179	(KD) aims to train a small student model to perform	Following previous work (Paetzold and Specia,	226
180	better by learning from a larger teacher model (Jian-	2016), we use precision and accuracy as metrics.	227
181	ping et al., 2021). We expect to learn a better small	<b>Precision (PRE):</b> The proportion of predicted al-	228
182	model from powerful LLMs, such as ChatGPT.	ternatives that are the original complex word itself	229
183	Since we can only access ChatGPT’s predictions,	or appear in the reference answers.	230
184	we adopt a black-box KD method that fine-tunes	<b>Accuracy (ACC):</b> The proportion of predicted	231
185	the student model on the data generated by the	alternatives that are different from the target com-	232
186	teacher model (Kim and Rush, 2016). To get high-	plex word and appear in the reference answers.	233
187	quality training data that are correct and diverse, we		
188	propose a pivot word-based approach for automatic	<b>3.3 Baseline Systems</b>	234
189	data generation based on ChatGPT.	We adopt BERT-LS (Qiang et al., 2021) and three	235
190	<b>Retrieval-augmented LLMs</b> OOD words are chal-	LLMs of various scales as the baselines. We try	236
191	lenging to simplify because the model may have	to analyze the behaviors and gain a deeper under-	237
192	little knowledge about them. Motivated by recent	standing of the advantages and disadvantages of	238
193	work on retrieval-augmented LLMs (Lewis et al.,	these models.	239
194	2020; Nakano et al., 2021), we propose a retrieval-		
195	based interpretation augmentation approach that	<b>3.3.1 BERT-LS</b>	240
196	dynamically brings in word interpretations from	The input of BERT-LS is formed by concatenat-	241
197	the web to enhance in-context learning.	ing the original sentence and its copy. The target	242
		complex word in the duplicate sentence is replaced	243
198	<b>3 Task, Data and In-Depth Analysis</b>	with [MASK]. BERT takes the input and predicts	244
		the masked part as substitution candidates.	245
199	In this section, we briefly introduce the task and the	Notice that a Chinese word often involves more	246
200	data, and analyze representative baselines which	than one Chinese character, while BERT’s tok-	247
201	are based on BERT and LLMs.	enizer is based on characters so BERT-LS allows	248
		BERT to make predictions through different num-	249
202	<b>3.1 Lexical Simplification Settings</b>	bers of [MASK] tags, e.g., from 1 to 4. All predic-	250
		tions are added to a list of candidates. When the tar-	251
203	An LS system first identifies complex words in a	get word is in a Chinese synonymy thesaurus (Mei,	252
204	sentence and then generates candidate substitutions,	1983), its synonyms are used as substitution can-	253
205	which is known as substitute generation (SG). Con-	didates. Finally, BERT-LS ranks these candidates	254
206	sidering complex word identification depends on a	with multiple sources of evidence including word	255
207	target population, we assume that a sentence and a	embeddings, BERT scores, and word frequencies.	256
208	target complex word are given following previous		
209	work (Qiang et al., 2021).	<b>3.3.2 LLMs</b>	257
210	Formally, given a sentence $s$ and a complex word	We use ChatGPT (GPT-4-1106-preview API),	258
211	$w$ in $s$ , the task is to generate a simpler alternative	ChatGLM2-6B (ChatGLM for short) (Du et al.,	259
212	$v$ , a word or a group of words, to form a simpler	2022) and ChatYuan-large-v2 (700m parameters,	260
213	sentence $s'$ , which is expected to be smooth, clear,	ChatYuan for short) (Xuanwei Zhang and Zhao,	261
214	and maintain the same meaning as $s$ .	2022). ChatGLM and ChatYuan are two open-	262
		source LLMs for dialogue, supporting both Chi-	263
215	<b>3.2 Dataset and Metrics</b>	nese and English languages.	264
		We explore the LLMs in a few-shot setting,	265
216	<b>3.2.1 Dataset</b>	where task instructions and demonstrations are in-	266
		cluded in the context. Figure 2 shows an exam-	267
217	We use the publicly available Chinese lexical sim-	ple. We use three demonstrations in this paper. For	268
218	plification dataset HanLS (Qiang et al., 2021).	ChatGPT, we extract predictions from its responses.	269
219	HanLS includes 524 sentences, each sentence con-	For ChatGLM and ChatYuan, we first extract the	270
220	tains a complex word from the advanced level		
221	of the Chinese Proficiency Test (Hanyu Shuiping		

Instruction	任务是将句子中给定的难词(由#标记)替换为一个简单的词或短语,同时保持句子的结构和意思不变并尽量流畅。 (The task is to replace the complex word (marked by #) in the sentence with a simple word or phrase, while keeping the structure and meaning of the sentence unchanged and as smooth as possible.)
Input 1	练习与理解不是#截然#对立的,而是相辅相成的。 (Practice and understanding are not #thoroughly# opposed but complement each other.)
Response	练习与理解不是#完全#对立的,而是相辅相成的。 (Practice and understanding are not #completely# opposed but complement each other.)
Input 2	他#呕心沥血#写了这本书。
Response	[Let LLM generate the response]

Figure 2: An example of instruction and demonstration design for prompting LLMs for CLS.

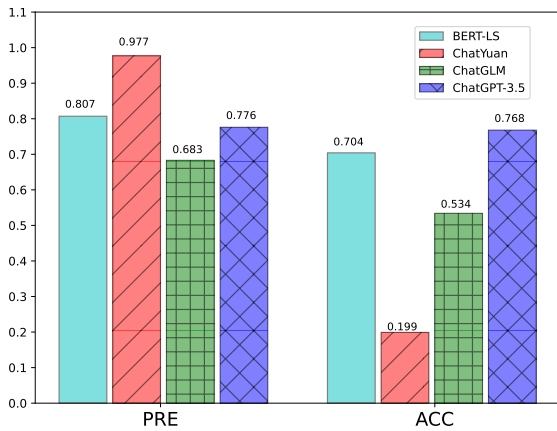


Figure 3: Overall results of BERT-LS and three LLMs in the few-shot setting.

top 10 candidates from model predictions and then re-rank them in the same way as BERT-LS.

## 3.4 Analysis and Discussion

### 3.4.1 Overall Results

Figure 3 shows the overall results of these systems. ChatYuan, BERT-LS, and ChatGLM are models that do not perform much simplification, as shown by their high PRE scores and low ACC scores. BERT-LS outperforms the other two models, possibly because for LS, MLM is a better fit for dialogue-oriented models in a few-shot setting. ChatGPT, on the other hand, simplifies almost all complex words and achieves high scores in both RPE and ACC, indicating its effectiveness.

Based on task instructions and demonstrations, ChatYuan and ChatGLM have difficulty in understanding the task. ChatGLM has a slight edge over ChatYuan, which may be attributed to its larger scale. The unsupervised BERT-LS also fails to

Models	Common		Idioms		OOD	
	PRE	ACC	PRE	ACC	PRE	ACC
BERT-LS	86.8	76.9	70.8	<b>41.7</b>	34.0	28.3
ChatYuan	<b>98.4</b>	22.4	<b>91.7</b>	8.30	<b>94.3</b>	3.80
ChatGLM	72.0	58.3	62.5	29.2	39.6	22.6
ChatGPT	81.6	<b>80.7</b>	29.2	29.2	66.0	<b>66.0</b>

Table 1: Detailed results of BERT-LS and three LLMs on simplifying 3 types of complex words.

Original sentence	小猪#似懂非懂#, 心想幸福怎么会是我的尾巴呢? The little pig #seemed to understand but didn't really understand#, thinking how could happiness be my tail?
BERT-LS	小猪#不解#, 心想幸福怎么会是我的尾巴呢? The little pig #puzzled#, thinking how could happiness be my tail?
ChatGPT	小猪#有些糊涂#, 心想幸福怎么会是我的尾巴呢? The little pig #was a bit confused#, thinking how could happiness be my tail?

Figure 4: The outputs of BERT-LS and ChatGPT on simplifying a Chinese idiom.

grasp the LS task, despite using external resources to re-rank candidates and compensate for its poor task understanding. ChatGPT, however, demonstrates a strong task understanding and a good performance in lexical simplification.

### 3.4.2 Analysis

We analyze the relation between the models' performance and the types of complex words. Specifically, we divide the complex words into 3 types:

- **Common words:** Refer to non-idiomatic words included in the dictionary Xinhua Zidian, which covers more than 320k words.
- **Chinese idioms:** Idioms or Chengyu, are an essential part of the Chinese language. They are usually composed of four Chinese characters and often express a moral or a lesson in a concise and elegant way.
- **Out-of-dictionary (OOD) words:** Refer to the words excluded in Xinhua Zidian, many of which are new words or internet terms.

Table 1 shows that ChatGPT outperforms BERT-LS, ChatYuan, and ChatGLM in simplifying common words, but lags behind BERT-LS on Chinese idioms. ChatGPT also has a remarkable advantage in simplifying OOD words, although none of the models achieve satisfactory results.

We compare the predictions of ChatGPT and BERT-LS in simplifying Chinese idioms and discover that ChatGPT's performance in simplifying

Chinese idioms is underrated because it generates phrases instead of single words as the reference answers do. For example, in Figure 4, ChatGPT produces a more understandable and fluent substitution than BERT-LS. BERT-LS only replaces the idiom with a single word, which may lose some descriptive details and degree of information.

Original sentence	我最近网上冲浪的时候总能刷到好多#镁铝#哦! I always see a lot of #magnesium aluminum# when I surf the internet recently!
BERT-LS	我最近网上冲浪的时候总能刷到好多#金属#哦! I always see a lot of #metal# when I surf the internet recently!
ChatGPT	我最近网上冲浪的时候总能刷到好多#美食#哦! I always see a lot of #gourmet food# when I surf the internet recently!

Figure 5: The outputs of BERT-LS and ChatGPT on simplifying an OOD word.

Simplifying OOD words is a challenge for all models. Figure 5 shows an example. The term “magnesium-aluminum” is a Chinese internet slang that sounds like “beauty” and refers to beautiful women. Neither BERT-LS nor ChatGPT can produce a suitable answer, probably because they have limited knowledge of these OOD words.

In summary, we discover the following observations through the in-depth analysis:

- **Task understanding:** Without enough supervision, the small and medium models, BERT-LS, ChatYuan, and ChatGLM, could not grasp the task well. ChatGPT shows a much better understanding of the task and performs well.
- **Sensitive performance to the type of complex words:** The difficulty of the simplification task depends on the types of complex words. BERT-LS and ChatGPT perform well in simplifying common complex words. For Chinese idiom simplification, ChatGPT has an advantage in preserving more descriptive details and degree of information. Simplifying OOD words is a challenge for all models.

ChatGPT performs the best but it is also expensive, which creates a trade-off between performance and cost. Moreover, it cannot deal with OOD words either. These factors lead us to the following research questions:

- **RQ 1:** How to use ChatGPT effectively as a teacher model to improve the task understanding and performance of smaller models on simplifying in-dictionary words?

- **RQ 2:** What are the effective strategies to enhance the performance of both large and small models in simplifying OOD words?
- **RQ 3:** What is the optimal way to integrate small and large models to achieve a trade-off between performance and cost?

## 4 The Proposed Method

We propose a framework as shown in Figure 1. It has 3 modules: automatic knowledge distillation from ChatGPT, retrieval-based interpretation augmentation, and a word type-aware controller.

### 4.1 Automatic Knowledge Distillation

We aim to create a high-quality CLS training dataset by distilling ChatGPT. We expect the generated sentences should be correct in spelling and grammar, cover diverse topics, and have accurate substitutions. However, since there are many factors to consider, it is unavoidable to bring in bias.

We propose an automatic knowledge distillation strategy named **PivotKD**, which only relies on ChatGPT and does not need any human intervention. Figure 6 illustrates its main workflow.

#### 4.1.1 Pivot Word Sampling

Our analysis shows that ChatGPT works well on common words and idioms but poorly on OOD words. So we should avoid OOD words in data generation. Therefore we sample words from the dictionary Xinhua Zidian and call the sampled words *pivot words*. We limit the word to be a noun, verb, adjective, adverb or idiom. To enhance diversity, each word can be sampled at most once. The pivot words would be used for pivot sentence generation.

#### 4.1.2 Pivot Sentence Generation

Given a pivot word, we ask ChatGPT to generate a sentence containing the word. The generated sentence would be used as the target sentence. This manner has the following benefits depending on the strengths of ChatGPT: 1) ChatGPT can generate correct and smooth sentences, avoiding spelling and grammar errors that often occur in data collected from the web or existing corpus; 2) ChatGPT can generate sentences covering diverse topics since we do not limit the topics in sampling pivot words and sentence generation, we can assume that the generated dataset is topic independent.

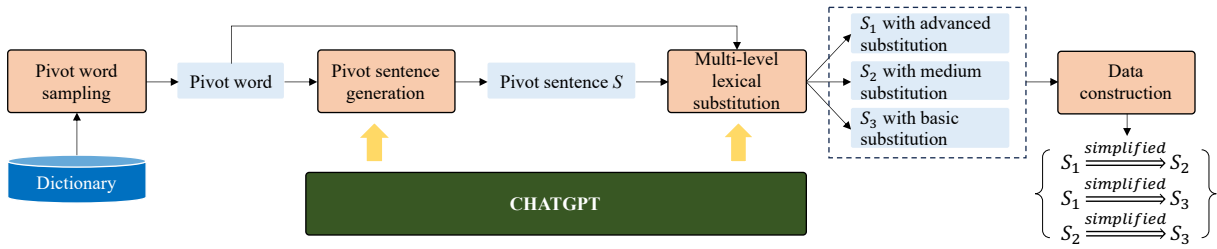


Figure 6: The main workflow of the PivotKD framework for generating CLS data based on ChatGPT.

任务是将句子中给定的难词(由#标记)替换为简单的词或短语,同时保持句子的结构和意思不变并尽量流畅。词语难度可以分为高级、中级和基本三个等级,请为每个等级分别生成 $n$ 个替换词。  
(We need to find a simpler word or phrase for the complex word, marked by #, in a sentence and keep the same structure and meaning of the original sentence while making the new sentence smooth. The word difficulty can be divided into three levels: advanced, intermediate, and basic. Please generate  $n$  replacement words for each level.)

Figure 7: An instruction for 3-level lexical substitution.

#### 4.1.3 Multi-level Lexical Substitution

After the above two steps, we have a pivot sentence and a pivot word. Now, we let ChatGPT generate substitutions belonging to three difficulty levels to replace the pivot word in the generated sentence. We define the three lexical difficulty levels as *advanced*, *medium*, and *basic*. We convey the requirements to ChatGPT through instructions as shown in Figure 7, and assume that ChatGPT itself can understand lexical difficulty.

#### 4.1.4 Data Construction

For one pivot word, we use ChatGPT to generate  $n$  simpler alternatives for each difficulty level. Given the  $3n$  sentences, we can construct a set of sentence pairs as training data based on the lexical difficulty levels of the substitutions. Specifically, a complex-to-simple sentence pair  $(s, s')$  can be constructed if the lexical difficulty level of the substitution in  $s$  is higher than the substitution in  $s'$ .

Notice that the sentence pair does not necessarily contain the pivot word. The pivot word only plays the role of starting point for the automatic knowledge distillation and data construction process, which implies the meaning of *pivot*.

#### 4.1.5 Instruction Fine-tuning

We conduct instruction fine-tuning with ChatYuan and ChatGLM. We fine-tune all parameters of ChatYuan and use LoRA (Hu et al., 2021) to fine-tune ChatGLM.

The training data is in the sequence-to-sequence style based on the constructed sentence pairs  $\{(s, s')\}$ . The input includes a task instruction, the same as the one shown in Figure 2, the target sentence  $s$ , and the target complex word marked with a tag # in  $s$ , while the output is the corresponding simplified sentence  $s'$  with the word substitution marked with # as well.

## 4.2 Retrieval-based Interpretation Augmentation

To ensure the quality of the generated data, we have focused on common words. However, simplifying OOD words is very challenging even for ChatGPT. To tackle this problem, we propose a retrieval-augmented strategy by looking for the interpretation of a target complex word from the web. **Retrieving Word Interpretation** Many OOD words are new words or internet slang that the pre-trained models may have limited knowledge of them. But there are usually interpretations for these words on the web. Therefore we utilize the Baidu search engine to fetch search results for the query "What does the word [complex word] mean?" and extract the content of the top  $k$  snippets as the interpretation.

**Injecting Interpretation for Inference** We use retrieval-based interpretation augmentation to provide additional context after the task instructions for ChatGPT or the fine-tuned models. So it is flexible and can be easily integrated as a plug-in.

## 4.3 Word-type aware Controlled Inference

Now we have several models and tools to handle CLS: ChatGPT, fine-tuned small models, and retrieval-based interpretation augmentation. We aim to find an effective and efficient way to integrate these modules.

Our solution is based on the observation that the performance of CLS is sensitive to the type of complex words. We prefer to use small models as the basic model and ask for help from ChatGPT

Models	Common		Idioms		OOD		All	
	PRE	ACC	PRE	ACC	PRE	ACC	PRE	ACC
BERT-LS	86.8	76.9	70.8	41.7	34.0	28.3	80.7	70.4
ChatGPT	81.6	80.7	29.2	29.2	66.0	66.0	77.6	76.8
+RIA	-	-	-	-	<b>79.3</b>	<b>79.3</b>	-	-
ChatGLM (frozen)	72.0	58.3	62.5	29.2	39.6	22.6	68.3	53.4
ChatGLM (fine-tuning with LoRA)	83.0	82.1	<b>66.7</b>	<b>66.7</b>	60.4	58.5	80.0	79.0
+RIA	<b>83.6</b>	<b>82.5</b>	58.3	58.3	<b>64.2</b>	<b>62.3</b>	<b>80.5</b>	<b>79.4</b>
ChatYuan (frozen)	98.4	22.4	91.7	8.3	94.3	3.8	97.7	19.9
ChatYuan (fine-tuning all parameters)	85.2	80.5	<b>75.0</b>	<b>70.8</b>	56.6	45.3	81.8	76.5
+RIA	<b>86.3</b>	<b>82.5</b>	<b>75.0</b>	<b>70.8</b>	<b>68.0</b>	<b>62.3</b>	<b>84.0</b>	<b>80.0</b>
Hyb-CLS	<b>86.3</b>	<b>82.5</b>	<b>75.0</b>	<b>70.8</b>	<b>79.3</b>	<b>79.3</b>	<b>85.1</b>	<b>81.6</b>

Table 2: System comparisons on HanLS. ChatGLM and ChatYuan are experimented with frozen, fine-tuning, and hybrid settings. RIA indicates utilizing retrieval-based interpretation augmentation during inference. The results with the highest accuracy are bolded, and the best results obtained by small models are marked with underlines.

or retrieval-based interpretation augmentation to handle Chinese idioms or OOD words. We design a rule-based word-type aware controller for deciding a proper inference strategy. We will discuss the optimal strategies for different types of words in the evaluation section.

## 5 Evaluation

### 5.1 Experimental Settings

For PivotKD, we sampled 5,000 pivot words from Xinhua Zidian. We avoid using the complex words in HanLS as pivot words. Since the sentences are fully generated by ChatGPT, there is also no overlap with HanLS.

The multi-level lexical substitution module generates  $n = 1$  substitution for each lexical difficulty level. Finally, we collect 8,962 sentence pairs, covering 4,269 distinct substitutions. More details about the data can be seen in Appendix A.

We conduct a human evaluation on 500 samples from the augmented dataset. For each complex word and simplified substitution pair, we let two persons judge whether the relative lexical difficulty between them are *clearly reasonable*, *hard to distinguish*, or *contradiction or irrelevant*. The proportions of the 3 options are 70%, 25%, and 5% respectively, indicating the quality of the dataset produced by PivotKD is acceptable.

We fine-tune ChatYuan for 1 epoch and ChatGLM with LoRA for 3 epochs. ChatYuan is based on T5, therefore the fine-tuned ChatYuan can be seen as the re-implementation of ConLS (Sheang et al., 2022). Detailed parameter settings are described in Appendix B. For retrieval-based interpretation augmentation, we use the top  $k = 1$  snippet

since it already obtains satisfactory performance.

## 5.2 Experimental Results

### 5.2.1 Auto-Evaluation

Table 2 shows the overall results and specific results on three types of complex words by BERT-LS, ChatGPT, ChatGLM, and ChatYuan and some variants. We can see some trends:

(1) **The effects of PivotKD** The fine-tuned ChatGLM and ChatYuan obtain competitive overall performance compared with ChatGPT and greatly outperform BERT-LS. For example, ChatGLM gets a 80.0 PRE score and a 79.0 ACC score, outperforming ChatGPT, while ChatYuan gets a 81.8 PRE score and a 76.5 ACC score.

The results demonstrate the effectiveness of PivotKD. The small models (e.g., ChatYuan) and medium models (e.g., ChatGLM) can benefit from supervised instruction fine-tuning based on the automatically generated data. The fine-tuned models understand the task much better and gain large improvements compared with the frozen models.

The effect of increasing the number of training samples on the performance of the fine-tuned ChatYuan and ChatGLM models is illustrated in Figure 8. Generally, the performance of the models can be enhanced by increasing the number of training samples. ChatGLM can reach a steady performance using relatively fewer samples, while the performance of ChatYuan has a consistent improvement with the increase of the training samples, indicating that smaller models may need more training data.

(2) **The effects of retrieval-based interpretation augmentation (RIA)** Table 2 shows that RIA can significantly enhance the ability of all three

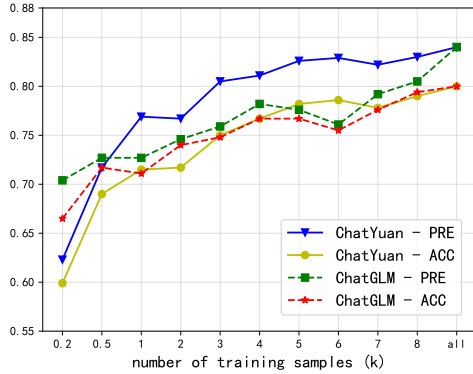


Figure 8: The effects of the number of training samples for fine-tuning ChatYuan and ChatGLM.

models to simplify OOD words, confirming that the retrieved word interpretation provides useful information for OOD words. ChatGLM and ChatYuan achieve similar performance to ChatGPT, while RIA boosts ChatGPT’s PRE and ACC scores by a large margin of 16%.

(3) **Performance on different word types** The performance of the fine-tuned models on different types of complex words is still very volatile, better on common words than Chinese idioms and OOD words. RIA also has different impacts on ChatGLM and ChatYuan. For example, ChatYuan gains further improvements for common words with RIA, while ChatGLM almost remains the same. Perhaps ChatGLM already has enough knowledge of the common words.

(4) **A hybrid approach** Since the performance is word-type sensitive, we can apply different settings for different types of words. We call the strategy **Hyb-CLS** (a hybrid approach for CLS). Specifically, we use the fine-tuned ChatYuan (700m) with RIA for common words and Chinese idioms, while call for ChatGPT to handle OOD words. Table 2 shows that the hybrid approaches obtain further improvements compared with any single model.

### 5.2.2 Human Evaluation

The system outputs may be reasonable but outside the reference answers. So we conduct a human evaluation. We sample 20 common words, 20 Chinese idioms, and 20 OOD words from HanLS. Three raters rate the mixed outputs of different systems according to the following criteria:

- 4 points: The substitution is simpler and has the same meaning as the target complex word without any information loss, and the resulting sentence is smooth.

Models	Common	Idioms	OOD	All
BERT-LS	3.17	1.23	0.6	1.67
ChatGPT + RIA	<b>3.70</b>	<b>3.17</b>	<b>2.60</b>	<b>3.16</b>
ChatGLM	3.37	2.60	1.93	2.63
+RIA	<u>3.50</u>	<u>2.83</u>	<u>2.53</u>	<u>2.95</u>
ChatYuan	3.37	2.60	1.50	2.49
+RIA	3.43	<u>2.83</u>	2.2	2.82

Table 3: Human evaluation of three models in different settings. The rating ranges from 0 (worst) to 4 (best).

- 2 points: The substitution is simpler and has the same meaning as the target word, but there is a loss of information in terms of details and degree, or the output is not so smooth.
- 0 points: The substitution is not simpler or its meaning is different from the target word.

Table 3 shows the averaged human evaluation results. We can see that ChatGPT still has an advantage in simplifying idioms and OOD words, indicating the strong ability of very large language models. The fine-tuned small models achieve similar performance and the performance is also close to ChatGPT. RIA is verified to be effective as well. The human evaluation confirms that with proper manipulation of the fine-tuned small models and very large models, it is possible to keep a balance between performance and cost.

## 6 Conclusion

This paper presents a word-type aware approach for Chinese lexical simplification. The core idea is to consider the types of complex words to effectively and efficiently combine small and large language models. We find that ChatGPT performs well in simplifying in-dictionary complex words and idioms. So we propose an automatic knowledge distillation framework called PivotKD to generate training data with ChatGPT for fine-tuning small models. The results show that the fine-tuned small models can outperform ChatGPT in simplifying such common words. Besides, we observe that both small and large models face challenges in simplifying OOD words. We propose a retrieval-based interpretation augmentation strategy, which significantly improves the simplification of OOD words for all models. Therefore, we can control the inference strategy according to the type of complex words, which efficiently combines small and large models and helps to obtain the best performance.



## 7 Limitations

There are three possible limitations of this work. First, our evaluation is based on the HanLS dataset, which is limited in size and coverage. We plan to extend the dataset. Second, we assume that ChatGPT understand the lexical difficulty levels, but we verify this assumption by analyzing the relative lexical difficulty between a pair of words in the generated data. More detailed and specially designed probing analysis can be conducted. Third, this paper focuses on Chinese lexical simplification, but the proposed method can be potentially applied to other languages. We plan to address these limitations in the future work.

## References

Dennis Aumiller and Michael Gertz. 2022. Unihd at tsar-2022 shared task: Is compute all we need for lexical simplification? In *Proceedings of the Workshop on Text Simplification, Accessibility, and Readability (TSAR-2022)*, pages 251–258.

John Carroll, Guido Minnen, Yvonne Canning, Siobhan Devlin, and John Tait. 1998. Practical simplification of english newspaper text to assist aphasic readers. In *Proceedings of the AAAI-98 Workshop on Integrating Artificial Intelligence and Assistive Technology*, pages 7–10. Citeseer.

Jan De Belder and Marie-Francine Moens. 2010. Text simplification for children. In *Proceedings of the SIGIR workshop on accessible search systems*, pages 19–26. ACM; New York.

Mohammad Dehghan, Dhruv Kumar, and Lukasz Golab. 2022. GRS: Combining generation and revision in unsupervised sentence simplification. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 949–960, Dublin, Ireland. Association for Computational Linguistics.

Biljana Drndarevic and Horacio Saggion. 2012. Towards automatic lexical simplification in spanish: An empirical study. In *PITR@ NAACL-HLT*, pages 8–16.

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.

Lijun Feng. 2009. Automatic readability assessment for people with intellectual disabilities. *ACM SIGACCESS accessibility and computing*, pages 84–91.

Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,

et al. 2021. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Gou Jianping, Yu Baosheng, Stephen J Maybank, and Tao Dacheng. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819.

Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.

Jiaju Mei. 1983. *Synonymy Thesaurus of Chinese Words*. Shanghai Lexicographical Publishing House, Shanghai.

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*.

Kai North, Marcos Zampieri, and Matthew Shardlow. 2023. Lexical complexity prediction: An overview. *ACM Computing Surveys*, 55(9):1–42.

OpenAI. 2023. Gpt-4 technical report.

Gustavo Paetzold and Lucia Specia. 2016. Unsupervised lexical simplification for non-native speakers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.

Gustavo Paetzold and Lucia Specia. 2017a. Lexical simplification with neural ranking. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 34–40.

Gustavo H Paetzold and Lucia Specia. 2017b. A survey on lexical simplification. *Journal of Artificial Intelligence Research*, 60:549–593.

Jipeng Qiang, Yun Li, Yi Zhu, Yunhao Yuan, and Xindong Wu. 2020. Lexical simplification with pre-trained encoders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8649–8656.

Jipeng Qiang, Xinyu Lu, Yun Li, Yunhao Yuan, and Xindong Wu. 2021. Chinese lexical simplification. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:1819–1828.

Horacio Saggion. 2017. Automatic text simplification. *Synthesis Lectures on Human Language Technologies*, 10(1):1–137.

Matthew Shardlow. 2013. A comparison of techniques to automatically identify complex words. In *51st Annual Meeting of the Association for Computational Linguistics Proceedings of the Student Research Workshop*, pages 103–109, Sofia, Bulgaria. Association for Computational Linguistics.

Kim Cheng Sheang, Daniel Ferrés, and Horacio Sag-gion. 2022. Controllable lexical simplification for english. In *Proceedings of the Workshop on Text Simplification, Accessibility, and Readability (TSAR-2022)*, pages 199–206.

Laura Vásquez-Rodríguez, Nhung Nguyen, Matthew Shardlow, and Sophia Ananiadou. 2022. Uom&mmu at tsar-2022 shared task: Prompt learning for lexical simplification. In *Proceedings of the Workshop on Text Simplification, Accessibility, and Readability (TSAR-2022)*, pages 218–224.

Liang Xu Xuanwei Zhang and Kangkang Zhao. 2022. Chatyuan: A large language model for dialogue in chinese and english. <https://github.com/clue-ai/ChatYuan>.

Seid Muhie Yimam, Chris Biemann, Shervin Malmasi, Gustavo Paetzold, Lucia Specia, Sanja Štajner, Anaïs Tack, and Marcos Zampieri. 2018. A report on the complex word identification shared task 2018. In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 66–78, New Orleans, Louisiana. Association for Computational Linguistics.

Attribute	Value
Sentence pairs	8,962
Avg. length of sentences	22.38
Distinct substitutions	4,269
Common words	4,186
Avg. length of substitutions	2.04

Table 4: Basic statistics of the augmented dataset via PivotKD.

## A Details in Dataset Construction

We sampled 5,000 pivot words for data generation. After constructing sentence pairs according to the difficulty levels of the substitutions, we use some rules to further reduce noise.

Firstly, we excluded substitutions that exist in the target complex word list of HanLS, thus there is no overlap between the augmented data and the test data. Secondly, we constrain that for the complex word in each constructed sentence pair should be in the Xinhua Zidian dictionary.

Some basic statistics of the final dataset for fine-tuning the small models are shown in Table 4. Table 5 shows a constructed sentence pairs and the corresponding training sample.

A sentence pair	
Complex	他#悄无声息#地走进房间，以免吵醒熟睡中的孩子。 (He slipped into the room #stealthily#, so as not to wake the sleeping child.)
Simple	他#偷偷#地走进房间，以免吵醒熟睡中的孩子。 (#quietly#.)
A training sample	
Prompt	他#悄无声息#地走进房间，以免吵醒熟睡中的孩子。 (Same as the sentence above)  你的任务是将句子中给定的难词#悄无声息#替换为一个简单的词或短语，同时保持句子结构和意思不变并尽量流畅。 (The task is to replace the complex word #stealthily# in the sentence with a simple word or phrase, while keeping the structure and meaning of the sentence unchanged and as smooth as possible. )
Response	他#偷偷#地走进房间，以免吵醒熟睡中的孩子。 (#quietly#.)

Table 5: An example of a constructed sentence pair and the corresponding training sample.

## B Parameter Settings

The hyper-parameters used for fine-tuning ChatYuan-large-v2 and ChatGLM2-6B are listed in Table 7 and Table 8 respectively.

We set the value of the temperature parameter of ChatGPT API as 0 because we emphasize the generation quality and control the diversity through pivot words.

## C Human Rating

We conducted human evaluation for ChatGLM, ChatYuan, ChatGPT and BERT-LS. We sampled 20 words for each type of complex word, merging the predictions of the models for human evaluation.

Settings	Value
GPU	Nvidia A6000
GPU memory	48 GB
CPU	AMD EPYC 7542
OS	Ubuntu 20.04.5 LTS
Pytorch version	1.31.1
CUDA version	11.6

Table 6: Infrastructure for conducting our experiments.

We set the ratings to be 0,2,or 4 according to the criteria introduced in the main content. We employed three raters who are students in a normal university. They are volunteers and unaware of the model information of these predictions. We reported the average rating for each prediction. The mean variance of the ratings between different raters is 0.55. Table 9 demonstrate two examples of predictions and human ratings.

Hyper-parameters	Value
max_seq_length (encoder)	512
max_seq_length (decoder)	512
num_epoch	1
learning_rate	5e-5
scheduler	cosine
batch_size	16
gradient_accumulation_steps	1

Table 7: Hyper-parameter settings used for fine-tuning ChatYuan-large-v2.

Hyper-Parameters	Value
max_seq_length (encoder)	512
max_seq_length (decoder)	512
num_epoch	3
learning_rate	5e-5
scheduler	cosine
batch_size	16
gradient_accumulation_steps	1
lora_rank	8
lora_alpha	32
lora_dropout	0.1

Table 8: Hyper-parameter settings used for fine-tuning ChatGLM2-6B.

Sentence	加强民族团结， #捍卫#国家完整。 (Enhance ethnic unity and #safeguard# national integrity.)
Prediction	加强民族团结， #维护(maintain)#国家完整。
Rater A	4
Rater B	4
Rater C	2
Sentence	那个年代，汤姆有一点叛逆， 有一个梦想就是去当 #绿林好汉#。 (In that era, Tom was a bit rebellious, and he had a dream of becoming an #outlaw hero#.)
Prediction	那个年代，汤姆有一点叛逆， 有一个梦想就是去当 #土匪(bandit)#。
Rater A	2
Rater B	0
Rater C	2

Table 9: Two examples of predictions and human ratings.