# Bilingual Terminology Extraction from Comparable E-Commerce Corpora

**Hao Jia**[1*], **Shuqin Gu**[2*], **Yuqi Zhang**[2], **Xiangyu Duan**[1‡]

[1]*Institute of Artificial Intelligence, School of Computer Science and Technology, Soochow University*
[2]*Machine Intelligence Technology Lab, Alibaba DAMO Academy*
hjia@stu.suda.edu.cn, shuqingu@tju.edu.cn, chenwei.zyq@alibaba-inc.com, xiangyuduan@suda.edu.cn

*Abstract*—Bilingual terminologies are important machine translation resources in the field of e-commerce, which are usually either manually translated or automatically extracted from parallel data. The human translation is costly and e-commerce parallel corpora is very scarce. However, the comparable data in different languages in the same commodity field is abundant. In this paper, we propose a novel framework of extracting e-commercial bilingual terminologies from comparable data. Benefiting from the cross-lingual pre-training in e-commerce, our framework can make full use of the deep semantic relationship between source-side terminology and target-side sentence to extract corresponding target terminology. Experimental results on various language pairs show that our approaches achieve significantly better performance than various strong baselines.

*Index Terms*—bilingual terminology extraction, e-commerce domain, cross-lingual pre-training

## I. INTRODUCTION

In recent years, many work has indicated that user-provided or domain-specific bilingual terminologies can enhance the accuracy and consistency of machine translation in specific domain [1]–[4]. Especially in the field of e-commerce [5], due to the diversity of product description, many terminologies[1] have their specific translations in specific product category, which makes it rather difficult for vanilla machine translation models to express their correct meanings. Moreover, the wrong translation of terminologies will lead to the decline of the whole sentence translation quality. As shown in Table I, the terminology "大款" (a terminology describing that the size of the cloth is big) is mistranslated into "big money" by Google Translate[2], and its correct translation in the sentence should be "big size". Obviously, when consumers browse this product on the e-commerce website, it will be misleading to consumers because of the wrong translation of the terminology describing the product attributes. Therefore, the correct translation of terminologies is of great significance to improving the translation quality in the field of e-commerce.

The acquisition of bilingual terminology pairs is either manual translation or automatically extracted from parallel data [6]–[10]. Manual translation is a reliable way, but it

---

* Hao Jia and Shuqin Gu make the equal contributions.

‡ Xiangyu Duan is the corresponding author.

[1]In this paper, e-commercial terminology refers to the key phrase that can describe product attributes, such as product name, product brand, product material, and product style, etc.

[2]https://translate.google.com/. We present the translation results on January 10th, 2022

---

TABLE I
EXAMPLE OF THE WRONG TERMINOLOGY TRANSLATION LEADING TO THE MISUNDERSTANDING OF E-COMMERCIAL SENTENCE.

| Source Sentence | 看来我只能买这种大款的 |
|---|---|
| Google Translate | It seems that i can only buy this kind of ***big money*** |
| Correct Translation | It seems that i can only buy this kind of ***big size*** |

is very time-consuming and expensive. The latter methods are either rule-based or statistical-based, using the linguistic feature, statistical feature or a hybrid of them. They rely on linguistic analysis tools, such as POS taggers, which may not be available for low-resource languages or domains.

The above automatic extraction methods are not suitable for e-commerce because of the lack of parallel e-commerce data. However, there are large-scale monolingual corpora covering different languages on popular e-commerce platforms. In such data, there are many potential terminology pairs, which are translations of each other. How to discover these bilingual terminologies is a big challenge in the e-commerce domain.

In this paper, we propose a new task, which is to discover bilingual terminologies from comparable data. The detailed description of constructing comparable data is presented in Section IV-A. Here, we focus on the e-commerce field. Given a terminology phrase in source language and a sentence in target language, the task is to 1) distinguish whether the target sentence contains the corresponding target translation of the source terminology, and 2) extract the corresponding target terminology from the target sentence if it contains.

To tackle this task, we propose an effective two-stage e-commercial bilingual terminology extraction framework. In the first stage, we fine-tune a cross-lingual pre-training model with a large number of e-commercial corpus consisting of different languages. In the second stage, we extract the target terminology from the target sentence by utilizing the extraction model initialized by cross-lingual pre-trained language models.

The main contributions of this paper can be summarized as follows:

- A new task of extracting bilingual terminologies of e-commerce from comparable data is proposed. In addition, we construct the corresponding comparable data in e-commerce domain.
- For the first time, the task of extracting bilingual terminologies of e-commerce from comparable data is for-

malized by using cross-lingual pre-training model and extraction framework.

- We conduct experiments mainly on three different e-commercial categories, namely clothes category, toys category and outdoors category in Chinese-to-English and English-to-French language pairs. Experimental results prove the effectiveness of the method. We hope our work would inspire new paradigms for bilingual terminology extraction.

## II. Related Work

### A. Cross-Lingual Word Embeddings for Bilingual Lexicon Induction

Following the success of word embeddings [11] trained on monolingual data, a large proportion of research aimed at mapping word embeddings into a common space for multiple languages [12]–[15], which were implemented by optimizing a linear transformation matrix. Based on these efforts, [16] proposed the extension of skip-gram to learn n-gram embeddings and mapped them to a shared space to obtain cross-lingual n-gram embeddings. However, these n-gram embeddings are based on the co-occurrence frequency.

### B. Bilingual Terminology Extraction from Parallel or Comparable Corpora

Several influential approaches [7], [9], [10], [17] have been proposed to extract bilingual terminology from parallel corpus, which mainly rely on the linguistic feature, statistical feature or the hybrid of them. [17] proposed an algorithm, which adopted English and French text taggers to associate noun phrases in the aligned English-to-French parallel corpus. The taggers provided part-of-speech categories which were used by finite-state recognizers to extract simple noun phrases for both languages. [18] proposed a sub-sentential alignment terminology extraction module that links linguistically motivated phrases in parallel texts. In addition, [19] proposed how to optimally combine different models derived from different resources for bilingual terminology extraction from comparable corpora. However, unlike our methods, these feature-driven (statistics or lingualistics) methods are usually not language-independent, and lack semantic information.

### C. Supervised Word Alignment Based on Cross-language Span Prediction

Researchers defined the alignment as an object for indicating the corresponding words in a parallel text [20], [21]. Recently, [22] formalized the supervised word alignment method as a cross-language span prediction problem similar to the SQuAD-style question answering task [23]. Specifically, given a target sentence as the context and a source word as a question, the word alignment system predicted a translation of the source word as the answer, which was a span in the target sentence.

Their idea is a little similar to our bilingual terminology extraction task based on cross-lingual pre-training model. However, in our method, in order to enhance the semantic



Fig. 1. Examples of our proposed task. Given a source-side terminology and a target-side sentence, we aim to distinguish whether the target-side sentence contains the translation of source-side terminology, and extract the corresponding translation if contains. Example 1 is the positive case, and Example 2 is the negative case.

representation, we utilize the e-commercial bilingual terminology pair and source term with corresponding target sentence pair to fine-tune the cross-lingual pre-training model. While [22] just utilized the multilingual BERT [24] as the semantic feature extractor.

## III. Proposed Task and Solutions

To our knowledge, we are the first to propose the task of e-commercial bilingual terminology extraction based on comparable corpus, independent on any parallel sentences. The definition of our proposed task and the solutions we proposed will be presented in the following.

### A. Task Definition

Our proposed task of e-commercial bilingual terminology extraction aims to extract the potential bilingual terminologies in massive e-commercial non-parallel corpus, which can be described as commonly specialized phrases with lengths of 2-5 grams.

In detail, given a terminology in source language and an e-commercial sentence in target language, we aim to distinguish whether the target sentence contains the corresponding target translation of the source terminology. If contains, we expect the system to find the position of the target terminology correctly. For example, as shown in Figure 1, the terminology in source language is "华为 P40 移动 5G手机", and the sentence in target language is "Global Version Huawei P40 Mobile 5G Phone 6.1 Inch Kirin 990 Android 10", the task is to predict the span of the potential terminology spans in the target sentence, i.e., "Huawei P40 Mobile 5G Phone", if not exists, return None.

### B. Our Approach

To tackle this task, we propose the two-stage e-commercial bilingual terminology extraction framework. In the first stage, we employ a large number of e-commercial corpus consists of different languages to perform cross-lingual pre-training. In the second stage, we extract the target terminology from the target sentence by utilizing Extractor_Attn or Extractor_Concat (illustrated in Figure 2) initialized by cross-lingual pre-trained language models in e-commerce.
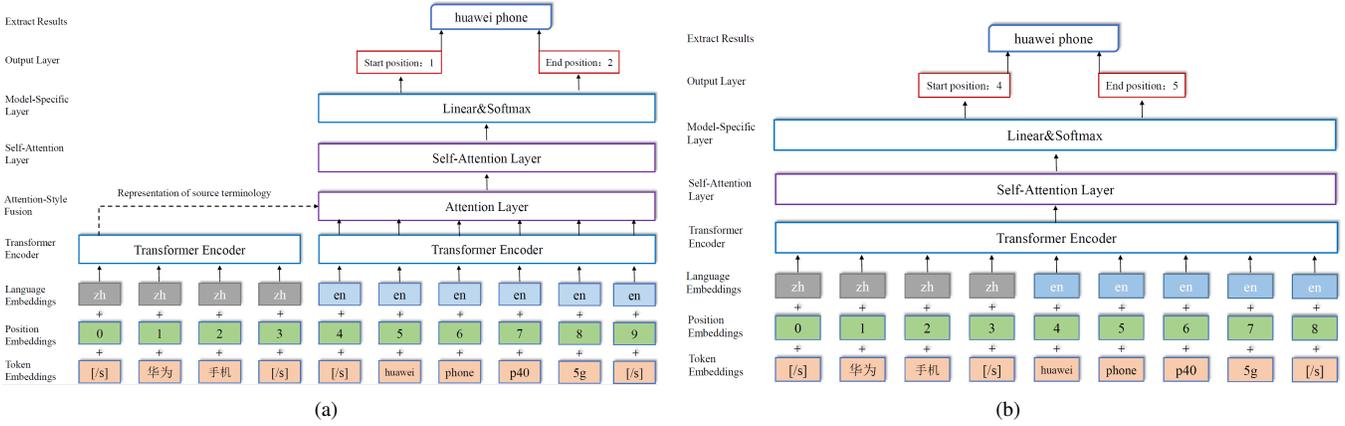
Fig. 2. The overview of our proposed methods, (a) Extractor_Attn and (b) Extractor_Concat.

## Cross-lingual Pre-training in E-Commerce

The cross-lingual language model pre-training (XLM) [25] method contains Masked Language Model (MLM) objective and Translation Language Model (TLM) objective, and has demonstrated its effectiveness on tasks such as XNLI cross-lingual classification and unsupervised machine translation. MLM is conducted over large monolingual corpora by randomly masking words, and training to predict them as a Cloze task [26]. Since MLM is only dependent on monolingual corpora, TLM is designed to utilize parallel data to drive better alignment between source and target language representations, by the means of concatenating parallel sentences into a whole sentence and randomly masking words in both the source and target side.

Inspired by these, we adopt MLM objective to perform cross-lingual pre-training on monolingual e-commercial corpus, which is the mixture of monolingual product titles in e-commerce domain from different languages. Besides, to gain better alignment between source and target language representations, we further propose to conduct TLM over the training sets of bilingual terminology pairs and the pairs of source terminology and target sentence in e-commerce.

## Target-side Terminology Extraction

Figure 2 generally illustrates our proposed framework. Given an e-commercial source terminology $S_{term}$ consisting of $m$ tokens $\{s_1, s_2...s_m\}$, we need to extract its corresponding translation span $T_{term}$ from the target sentence $T$ containing $n$ tokens $\{t_1, t_2...t_n\}$ . We use the Transformer [27] encoder initialized by cross-lingual pre-trained models in e-commerce as the backbone to fully extract the deep semantic relationship between the source-side terminology and target-side sentence, so that our framework could correctly distinguish or even extract the target-side terminology. We propose two methods Extractor_Attn and Extractor_Concat to proceed the extraction of representation.

*Extractor_Attn*: As illustrated in Figure 2(a), $S_{term} = \{[/s], s_1, s_2...s_m, [/s]\}$ and $T = \{[/s], t_1, t_2...t_n, [/s]\}$, consisting of the adding sum of language embedding, position

embedding and token embedding of corresponding tokens, are fed into the Transformer encoder respectively to get the representation matrix $H_{src\_term} \in \mathbb{R}^{(m+2) \times d}$ and $H_{tgt\_stc} \in \mathbb{R}^{(n+2) \times d}$. Then we obtain the final representation matrix $H$ by doing representational fusion between source-side terminology and target-side sentence as follows:

$$F(H_{tgt\_stc}) = MultiHead(Q = H_{tgt\_stc}, \\ K = H_{src\_term}, V = H_{src\_term}) \qquad (1)$$

$$H_1 = FFN(LayerNorm(H_{tgt\_stc} + F(H_{tgt\_stc}))) \qquad (2)$$

$$F(H_1) = MultiHead(Q = H_1, K = H_1, V = H_1) \qquad (3)$$

$$H = FFN(LayerNorm(H_1 + F(H_1))) \qquad (4)$$

where $MultiHead$, $LayerNorm$, $FFN$ are basic components of the Transformer model. By the attention-style fusion and self-attention computation in this way, the model can utilize the representation of source-side terminology as weight to attend the most related span of target-side sentence. Note that the encoders of source-side terminology and target-side sentence share the same parameters.

*Extractor_Concat*: As illustrated in Figure 2(b), the input sequence consists of $S_{term}$ concatenated with $T$, i.e., $\{[/s], s_1, s_2...s_m, [/s], t_1, t_2...t_n, [/s]\}$ where $[/s]$ is a special token. Then the Transformer encoder utilizes the embedding representation of input sequence, which is calculated as the adding of language embedding, position embedding and token embedding of corresponding tokens, to perform attention-style fusion and self-attention computation. As a result, the encoder will outputs a cross-lingual context representation matrix $H \in \mathbb{R}^{(m+n+3) \times d}$, where $d$ is the vector dimension of the last layer of the encoder. In this way, the model can attend to both source-side terminology and target-side sentence, encouraging the model to learn and align the source and target representations.

**Span Detector:** Given the representation matrix output $H$ from Extractor_Attn/Extractor_Concat, we then input it to the linear layer, so as to separately predict the start index and the

end index of the target terminology in target sentence. It can be formulated as follows:

$$\mathbf{p}_{start} = softmax(\mathbf{W}_{start} \cdot H) \tag{5}$$

$$\mathbf{p}_{end} = softmax(\mathbf{W}_{end} \cdot H) \tag{6}$$

where both of the $\mathbf{W}_{start} \in \mathbb{R}^{d \times 2}$ and $\mathbf{W}_{end} \in \mathbb{R}^{d \times 2}$ are linear layers with learnable parameters.

**Loss Function:** During the training, we separately calculate the loss of predicting the start index and end index of the target terminology, which are given as follows:

$$\mathcal{L}_{start} = \mathbf{CE}(\mathbf{p}_{start}, \mathbf{y}_{start}) \tag{7}$$

$$\mathcal{L}_{end} = \mathbf{CE}(\mathbf{p}_{end}, \mathbf{y}_{end}) \tag{8}$$

where **CE(\*)** refers to cross-entropy computation. Then, the overall training objective to be minimized is as follows:

$$\mathcal{L} = \frac{1}{2}(\mathcal{L}_{start} + \mathcal{L}_{end}) \tag{9}$$

The two losses are jointly trained, with parameters shared at the linear layer.

Note that in the **inference** phase, the start and end indexes will be predicted respectively. If both of them equal 0 or the start index is larger than the end index, it means that there are no corresponding target terminology in the current sentence. If not, leading to the final extracted results of target terminology.

## IV. EXPERIMENTS

We conduct experiments on Chinese→English and English→French e-commercial corpus to demonstrate the effectiveness of our proposed solutions.

### A. Data Construction

In this section, we describe the process of constructing e-commercial comparable corpus in detail. We acquire English, Chinese, and French monolingual texts from the popular e-commerce platforms covering three product categories: clothes, toys, and outdoors. Monolingual texts under the same product category in different languages can be seen as e-commercial comparable corpus, which is not parallel sentence pairs but may contain potential parallel terminology pairs.

For each product category in all language pairs, we select frequent e-commercial terminologies from the monolingual sentences in source language (i.e., Chinese and English), and manually translate them into target language (i.e., English and French), which constitutes bilingual terminologies. In addition, we retrieve the monolingual sentences in target language containing the target terminology. If the target sentence contains the target terminology, we can construct a data pair of source terminology, target terminology and target sentence, which can be noted as the positive case, if not, negative case instead.

TABLE II
STATISTICS OF THE DATE SETS.

| Data Sets | | E-commercial Categories | | |
|---|---|---|---|---|
| | | clothes | toys | outdoors |
| zh→en | training set | 0.68M | 0.46M | 0.60M |
| | validation set | 1000 | 1000 | 1000 |
| | test set | 2694 | 2426 | 2396 |
| en→fr | training set | 0.61M | 0.61M | 0.61M |
| | validation set | 1000 | 1000 | 1000 |
| | test set | 2266 | 2442 | 2334 |

### B. Experimental Setup

*a) Data Sets:* Following the data construction method described in Section IV-A, we construct data of e-commercial bilingual terminology pairs, and the positive and negative cases of {*source terminology, target terminology, target sentence*} pairs. For positive cases, we get the start and end indices of the target terminology in corresponding target sentence. For negative cases, we set both the start and end indices as 0. Sequentially, we divide these data pairs into training, validation and test sets. In training/validation/test sets, the ratio of positive and negative cases remains at 1:1. The statistics of data sets are summarized in Table II.

For cross-lingual pre-training in e-commerce, we use all the available monolingual e-commercial title corpus to perform MLM, which contains 5.5M, 7.4M, 4.1M for English, Chinese and French, respectively. Specially, we conduct TLM over the bilingual terminology pairs and the positive portion of the training set. The training sets of bilingual terminology pairs consist of 21,500 for Chinese→English and 19,500 for English→French.

*b) Training Configuration:* For cross-lingual pre-training stage, we conduct joint byte-pair encoding (BPE) on the monolingual or comparable corpora of both languages with a shared vocabulary. We use the cross-lingual pretrained models released by XLM[3] for the model initialization. During pre-training, following Conneau and Lample [25], 15% of BPE tokens are selected to be masked. Among the selected tokens, 80% of them are replaced with [MASK] token, 10% are replaced with a random BPE token within the vocabulary, and 10% remain unchanged.

For target terminology extraction phase, we adopt the commonly used Transformer encoder with 1024 embedding/hidden units, 4096 feed-forward filter size, 6 layers and 8 heads per layer as the basic. During training, the batch size is set to 128 and the sentence length is limited to 100 BPE tokens. We employ the Adam [28] optimizer with $lr = 0.0001$, $t_{warm\_up} = 4000$ and $dropout = 0.1$.

*c) Evaluation Metric:* During evaluating, we calculate the accuracy whether the model correctly predict both the start and end indices of the target-side terminology as follows:

$$accuracy = \frac{Nums_{correct}}{Nums_{all}} \tag{10}$$

[3]https://github.com/facebookresearch/XLM

| System | | zh→en | | | | en→fr | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | clothes | toys | outdoors | avg. | clothes | toys | outdoors | avg. |
| Baselines | NMT | 51.51 | 44.42 | 43.32 | 46.42 | 31.60 | 28.09 | 28.00 | 29.23 |
| | SMT | 54.31 | 47.13 | 47.63 | 49.69 | 39.81 | 36.69 | 36.82 | 37.77 |
| | Multiple MT Voting | 47.92 | 56.22 | 53.84 | 52.66 | 53.40 | 57.41 | 56.93 | 55.91 |
| | Seq2Seq-Term | 65.26 | 50.78 | 54.01 | 56.68 | 41.66 | 30.06 | 32.65 | 34.79 |
| | CLSP$_{eco}$ [22] | 86.06 | 73.52 | 74.81 | 78.13 | 70.84 | 64.92 | 68.46 | 68.07 |
| Extractor_Attn | RAND | 77.13 | 56.55 | 57.51 | 63.73 | 53.84 | 41.69 | 43.19 | 46.24 |
| | MLM$_{eco}$ | 84.86 | 70.57 | 73.54 | 76.32 | 67.87 | 62.33 | 61.30 | 63.83 |
| | TLM$_{eco}$ | 85.67 | 74.11 | 75.38 | 78.39 | 70.52 | 66.50 | 70.35 | 69.12 |
| Extractor_Concat | RAND | 83.96 | 69.41 | 70.45 | 74.61 | 76.00 | 62.16 | 66.36 | 68.17 |
| | MLM$_{eco}$ | 92.92 | 87.07 | 86.48 | 89.59 | 90.86 | 88.92 | 87.05 | 88.94 |
| | TLM$_{eco}$ | **94.43** | **91.51** | **90.23** | **92.06** | **92.58** | **90.25** | **90.54** | **91.12** |

, where $Nums_{correct}$ denotes the number of cases correctly predicted by the model, and $Nums_{all}$ denotes the number of all cases in the test set.

### C. Baselines

We adopt the following methods as our baselines:

- **NMT/SMT:** We take the task of bilingual terminology extraction as an MT problem, i.e., bilingual terminology generation. Source terminology is fed into the MT model and the output sequence is target terminology. We adopt Transformer[4] [27] and Moses[5] as the NMT model and SMT system respectively. We measure whether the model correctly generates the entire target terminology as equation 10.
- **Multiple MT Voting:** We firstly utilize *Google*[6], *Baidu*[7], *Youdao*[8], *bing*[9] and *sogou*[10] Translate Systems to directly translate the source terminology in our test set and get the corresponding translation candidates. Then we vote according to the results of different MT systems, with the highest number of votes as the final translation. We measure whether the final translation is the correct target terminology as equation 10.
- **Cross-Language Span Prediction in E-Commerce (CLSP$_{eco}$) [22]:** Cross-language span prediction method has been used for neural word alignment [22], which can also be applied in e-commercial bilingual terminology extraction. To adapt it to e-commerce domain, we fine-tune multilingual BERT with e-commercial monolingual corpora. Then we formalize the task as SQuAD-style span prediction problem and solve it with the fine-tuned multilingual BERT as they propose in the paper.

- **Seq2Seq-Term:** We regard the task as a sequence-to-sequence (seq2seq) learning problem by encoding the input of source terminology concatenated with target sentence, and decoding the output sequence of target terminology. For positive cases, the model will decode the corresponding target terminology. While for negative cases, the model will decode a special token "None", which means the translation of the source terminology does not exist in target sentence. Our Seq2Seq-Term system builds on Transformer[11], a state-of-the-art seq2seq model, with the shared vocabulary between input and output. This baseline is most related to our approaches, since they utilize the same data resources.

### D. Main Results

Table III presents the performance of our proposed approach and other baseline models on different categories of different language pairs. It is obvious that our approaches outperform the baselines significantly in all language pairs and categories, which strongly demonstrates the superiority of cross-lingual pre-training and our proposed bilingual terminology extraction models.

**Comparison between Baselines**

The performances of SMT systems are consistently superior to NMT systems, which indicates that directly using SMT trained on bilingual terminology pairs is more suitable for the task of bilingual terminology generation than NMT. In particular, we can find that Multiple MT Voting achieves better performance, mainly because it acquires the final translation results by voting among several top-tier MT engines. Seq2Seq-Term performs best among all baselines in zh→en, while worse than Multiple MT Voting and SMT in en→fr.

**Comparison between Baselines and Our Proposed Approaches**

---

[4]https://github.com/pytorch/fairseq/tree/v0.9.0. We use Transformer$_{base}$ as our model.

[5]http://www.statmt.org/moses/. We use the default setting of Moses.

[6]https://translate.google.com/

[7]https://fanyi.baidu.com/

[8]http://fanyi.youdao.com/

[9]https://bing.com/translator

[10]https://fanyi.sogou.com/text

[11]https://github.com/pytorch/fairseq/tree/v0.9.0. We use Transformer$_{base}$ as our model.

| Source terminology | Target sentence | Polarity |
|---|---|---|
| 两@@ 件 套装 | [/s] men@@ s tr@@ ack@@ suit sets brand two piece suit tr@@ ack@@ suit 2019 male cas@@ ualt@@ shir@@ ts [/s] | Positive |
| 儿童 衬@@ 衣 | [/s] al@@ phal@@ moda 2019 new flor@@ al cot@@ ton dress v-@@ neck ru@@ ffled sle@@ eve high waist a-@@ line [/s] | Negative |

Fig. 3. Visualized attention weights for source-side terminology and target-side sentence by Extractor_Concat. "Positive" and "Negative" in column "Polarity" indicate whether the target sentence contains the corresponding translation of the source terminology or not.

Compared with various baselines, our proposed Extractor_Attn and Extractor_Concat with random initialization both show clear superiority, which demonstrates the effectiveness of our proposed bilingual terminology extraction models. Especially in comparison with the most related Seq2Seq-Term, our models show better performances, indicating that Extractor_Attn and Extractor_Concat are more suitable for the task of bilingual terminology extraction than the seq2seq method.

**Comparison among Different Initialization Methods**

When equipped with $MLM_{eco}$ or $TLM_{eco}$, our proposed Extractor_Attn and Extractor_Concat gain great improvements (+ 7.73%-28.09%), proving the significance of cross-lingual pre-training in e-commerce for the extraction models. Specially, models initialized with $TLM_{eco}$ perform consistently better than those initialized with $MLM_{eco}$ in all product categories. Obviously, the models could learn rich cross-lingual alignment information by $TLM_{eco}$, which encourages the extraction models to better distinguish and even extract the target-side terminology.

**Comparison between Extractor_Attn and Extractor_Concat**

Particularly, when comparing Extractor_Attn and Extractor_Concat, we note that Extractor_Concat outperforms Extractor_Attn under all model initialization conditions. Moreover, Extractor_Concat initialized with $TLM_{eco}$ obtains the best performances in all languages and all categories. It is because that Extractor_Concat conducts self-attention computation on both source-side terminology and target-side sentence at the same time, while Extractor_Attn calculates self-attention on source-side terminology and target-side sentence separately. We argue that Extractor_Concat learns richer cross-lingual semantic relationship between source terminology and target sentence, and pay more attention to the most related span of target-side sentence.

TABLE IV
PERFORMANCES(%) OF EXTRACTOR_CONCAT WITH OR WITHOUT SOURCE-SIDE TERMINOLOGIES, WITH DIFFERENT INITIALIZATION PARAMETERS IN CHINESE→ENGLISH.

| | | clothes | toys | outdoors | avg. |
|---|---|---|---|---|---|
| w/o source term | RAND | 31.55 | 31.90 | 38.81 | 34.09 |
| | $MLM_{eco}$ | 35.12 | 33.97 | 40.48 | 36.52 |
| | $TLM_{eco}$ | 35.78 | 34.38 | 40.90 | 37.02 |
| w/ source term | RAND | 83.96 | 69.41 | 70.45 | 74.61 |
| | $MLM_{eco}$ | 92.92 | 87.07 | 86.48 | 89.59 |
| | $TLM_{eco}$ | 94.43 | 91.51 | 90.23 | 92.06 |

TABLE V
PERFORMANCES(%) OF EXTRACTOR_CONCAT WITH OR WITHOUT THE TOP SELF-ATTENTION LAYER, WITH DIFFERENT INITIALIZATION PARAMETERS IN CHINESE→ENGLISH.

| | | clothes | toys | outdoors | avg. |
|---|---|---|---|---|---|
| w/ self-attn layer | RAND | 83.96 | 69.41 | 70.45 | 74.61 |
| | $MLM_{eco}$ | 92.92 | 87.07 | 86.48 | 89.59 |
| | $TLM_{eco}$ | 94.43 | 91.51 | 90.23 | 92.06 |
| w/o self-attn layer | RAND | 83.67 | 68.26 | 69.78 | 73.90 |
| | $MLM_{eco}$ | 92.72 | 86.69 | 85.64 | 88.35 |
| | $TLM_{eco}$ | 94.21 | 91.18 | 90.12 | 91.84 |

## V. ANALYSIS

### A. Effect of Source Terminology

In our proposed Extractor_Concat, source terminology and target sentence are concatenated as an input sequence to the model, forming the final representation after self-attention computation. We wonder whether the model really learns the semantic relationship between source terminology and target sentence, or just extracts the target terminology depending on target sentence by simple positional recognition. Therefore, we remove the source terminology and take only the target sentence as input, attempting to predict target terminology just dependent on target sentence. Table IV shows the performance of Extractor_Concat with or without source-side terminologies in Chinese→English. We can observe that without source-side terminologies, the performances drop notably, which

demonstrates the effect of interactiveness between source-side terminologies and target-side sentences.

### B. Effect of Last Self-attention Layer

In our proposed Extractor_Concat, self-attention layer is employed on the encoder output to obtain the cross-lingual context representation. We do ablation study to show the contribution of the top self-attention layer. Table V the performance of Extractor_Concat with or without the top self-attention layer in Chinese→English. It shows that the performance of Extractor_Concat decreases when the top self-attention layer is removed, which demonstrates the significance of self-attention layer.

### C. Contextualized Word Representation

To further investigate the effects of cross-lingual alignment, we sample two pairs of source terminology and target sentence from the Chinese→English validation sets, and visualize the attention weights on them in Figure 3. The color depth indicates the importance degree of the weight, the darker the more important. As can be seen, the semantic similarity between source and target terminology are able to be captured. In the positive example, 两@@ 件 套装 matches *two piece suit*, which are mutual translations. While in the negative example, 儿童 衬@@ 衣 matches *[/s]* in the beginning of target sentence, since target sentence does not contain its translation in target language.

### D. Effect of Training Data Size

We expand the number of positive and negative cases in the training sets (the ratio remains 1:1), so as to study the effect of the training data size on the performance. Table 4 reports the performances of our proposed Extractor_Concat with double, triple and quadruple training data, i.e., 1.2M, 1.8M or 2.4M. It shows that when we expand the size of training data to twice, the performances of all categories improve significantly. But when the size of training data is expanded to triple and quadruple, the performances of all categories drop sharply. More training data will not always lead to improvements of systems performances, suggesting that our model have already learnt enough cross-lingual semantic information with limited training data.

### VI. Conclusion

In this paper, we propose a new task of extracting bilingual terminologies from non-parallel comparable corpus in e-commerce and construct corresponding data sets. We apply a two-stage neural framework to tackle this task. When equipped with cross-lingual pre-training in e-commerce, our proposed Extractor_Concat and Extractor_Attn can extract the corresponding target terminology by fully utilizing the deep semantic relationship between source-side terminology and target-side sentence. Experimental results show that our methods outperform all strong baselines in all categories on the Chinese→English and English→French language pairs. As far as we know, we are the first to utilize cross-lingual pre-training
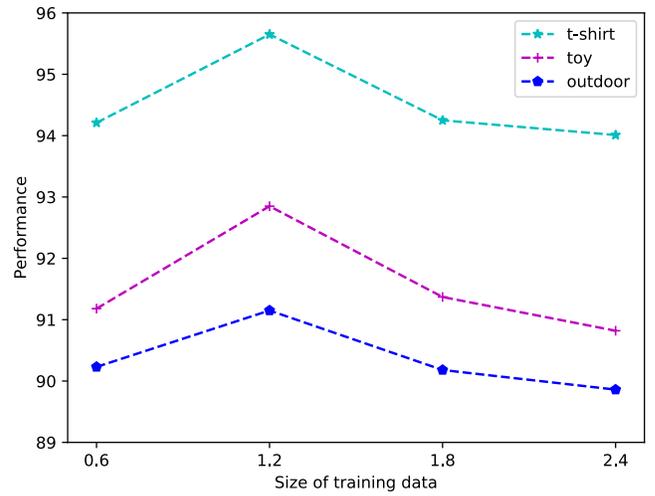


Fig. 4. Performances(%) of varying size(M) of training samples for Extractor_Concat initialized by TLM$_{eco}$.

and extraction model to solve the problem of extracting bilingual terminologies from non-parallel e-commerce corpora. We hope that our work will encourage the introduction of new paradigms for bilingual terminology extraction or other relevant research.

### References

[1] E. Hasler, A. De Gispert, G. Iglesias, and B. Byrne. 2018. Neural machine translation decoding with terminology constraints. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 506-512.

[2] C. Hokamp, and Q. Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1535-1546.

[3] M. Post, and D. Vilar. 2018. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1314-1324.

[4] P. Arthur, G. Neubig, and S. Nakamura. 2016. Incorporating discrete translation lexicons into neural machine translation. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1557-1567.

[5] K. Song, Y. Zhang, H. Yu, W. Luo, K. Wang, and M. Zhang. 2019. Code-switching for enhancing nmt with pre-specified translation. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 449-459.

[6] É. Gaussier. 1998. Flow network models for word alignment and terminology extraction from bilingual corpora. In COLING 1998 Volume 1: The 17th International Conference on Computational Linguistics.

[7] É. Gaussier, D. Hull, and S. Ait-Mokhtar. 2000. Parallel Text Processing Alignment and Use of Translation Corpora, chapter Term Alignment in Use: Machine-Aided Human Translation.

[8] D. Chambers. 2000. Automatic Bilingual Terminology Extraction: A Practical Approach. In Proceedings of Translating and the Computer, 22.

[9] G. F. Le An Ha, R. Mitkov, and G. Corpas. 2008. Mutual, bilingual terminology extraction. In Proceedings of LREC 2008, Marrakesh. Citeseer.

[10] R. Haque, S. Penkale, and A. Way. 2014. Bilingual termbank creation via log-likelihood comparison and phrase-based statistical machine translation. In Proceedings of the 4th International Workshop on Computational Terminology (Computerm), 42-51.

[11] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. 2013. Distributed representations of words and phrases and their compositionality. Advances in neural information processing systems, 26: 3111-3119.

[12] M. Zhang, Y. Liu, H. Luan, and M. Sun. 2017. Adversarial training for unsupervised bilingual lexicon induction. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 1959-1970.

[13] A. Conneau, G. Lample, M. Ranzato, L. Denoyer, and H. Jégou. 2017. Word Translation Without Parallel Data. arXiv preprint arXiv:1710.04087.

[14] G. Lample, A. Conneau, L. Denoyer, and M. Ranzato. 2017. Unsupervised Machine Translation Using Monolingual Corpora Only. arXiv preprint arXiv:1711.00043.

[15] A. Lazaridou, G. Dinu, and M. Baroni. 2015. Hubness and pollution: Delving into cross-space mapping for zero-shot learning. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), 270-280.

[16] M. Artetxe, G. Labaka, and E. Agirre. 2018. Unsupervised Statistical Machine Translation. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 3632-3642.

[17] J. Kupiec. 1993. An algorithm for finding noun phrase correspondences in bilingual corpora. In 31st Annual Meeting of the Association for Computational Linguistics, 17-22.

[18] E. Lefever, L. Macken, and V. Hoste. 2009. Language-independent bilingual terminology extraction from a multilingual parallel corpus. In Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009), 496-504.

[19] H. Déjean, É. Gaussier, and F. Sadat. 2002. Bilingual terminology extraction: an approach based on a multilingual thesaurus applicable to comparable corpora. In Proceedings of the 19th International Conference on Computational Linguistics COLING, 218-224. Citeseer.

[20] F. J. Och, and H. Ney. 2003. A systematic comparison of various statistical alignment models. Computational linguistics, 29(1): 19-51.

[21] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions, 177-180.

[22] M. Nagata, C. Katsuki, and M. Nishino. 2020. A supervised word alignment method based on cross-language span prediction using multilingual BERT. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 555-565.

[23] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. arXiv preprint arXiv:1606.05250.

[24] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. 2019. 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), pages 4171-4186.

[25] A. Conneau, and G. Lample. 2019. Cross-lingual language model pretraining. In Advances in Neural Information Processing Systems, 7059-7069.

[26] W. L. Taylor. 1953. Cloze procedure: A new tool for measuring readability. Journalism quarterly, 30(4): 415-433.

[27] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. 2017. Attention is All you Need. In Advances in Neural Information Processing Systems.

[28] D. P. Kingma, and J. Ba. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.