



Original Article

Adaptation in Statistical Machine Translation for Low-resource Domains in English-Vietnamese Language

Nghia-Luan Pham^{1,2,*}, Van-Vinh Nguyen²

¹*Hai Phong University, 171 Phan Dang Luu, Kien An, Haiphong, Vietnam*

²*Faculty of Information Technology, VNU University of Engineering and Technology,
Vietnam National University, Hanoi, 144 Xuan Thuy, Cau Giay, Hanoi, Vietnam*

Received 09 April 2019

Revised 19 May 2019; Accepted 13 December 2019

Abstract: In this paper, we propose a new method for domain adaptation in Statistical Machine Translation for low-resource domains in English-Vietnamese language. Specifically, our method only uses monolingual data to adapt the translation phrase-table, our system brings improvements over the SMT baseline system. We propose two steps to improve the quality of SMT system: (i) classify phrases on the target side of the translation phrase-table use the probability classifier model, and (ii) adapt to the phrase-table translation by recomputing the direct translation probability of phrases.

Our experiments are conducted with translation direction from English to Vietnamese on two very different domains that are legal domain (*out-of-domain*) and general domain (*in-of-domain*). The English-Vietnamese parallel corpus is provided by the IWSLT 2015 organizers and the experimental results showed that our method significantly outperformed the baseline system. Our system improved on the quality of machine translation in the legal domain up to 0.9 BLEU scores over the baseline system,...

Keywords: Machine Translation, Statistical Machine Translation, Domain Adaptation.

1. Introduction

Statistical Machine Translation (SMT) systems [1] are usually trained on large amounts of bilingual data and monolingual target language data. In general, these corpora

may include quite heterogeneous topics and these topics usually define a set of terminological lexicons. Terminologies need to be translated taking into account the semantic context in which they appear.

The Neural Machine Translation (NMT) approach [2] has recently been proposed for machine translation. However, the NMT method requires a large amount of parallel data and it has some characteristics such as NMT

* Corresponding author.

E-mail address: luanpn@dhhp.edu.vn

<https://doi.org/10.25073/2588-1086/vnucsce.231>

system is too computationally costly and resource, the NMT system also requires much more training time than SMT system [3]. Therefore, SMT systems are still being studied for specific domains in low-resource language pairs.

Monolingual data are usually available in large amounts, parallel data are low-resource for most language pairs. Collecting sufficiently large-high-quality parallel data is hard, especially on domain-specific data. For this reason, most languages in the world are low-resource for statistical machine translation, including the English-Vietnamese language pair.

When SMT system is trained on the small amount of specific domain data leading to narrow lexical coverage which again results in low translation quality. On the other hand, the SMT systems are trained, tuned on specific-domain data will perform well on the corresponding domains, but performance deteriorates for out-of-domain sentences [4].

Therefore, SMT systems often suffer from domain adaptation problems in practical applications. When the test data and the training data come from the same domains, the SMT systems can achieve good quality. Otherwise, the translation quality degrades dramatically. Therefore, domain adaptation is of significant importance to developing translation systems which can be effectively transferred from one domain to another.

In recent years, the domain adaptation problem in SMT becomes more important [5] and is an active field of research in SMT with more and more techniques being proposed and applied into practice [5-12]. The common techniques used to adapt two main components of contemporary state-of-the-art SMT systems: The language model and the translation model. In addition, there are also some proposals for adapting the Neural Machine Translation (NMT) system to a new domain [13, 14]. Although the NMT system has begun to be studied more, domain adaptation for the SMT system still plays an important role, especially for low-resource languages.

This paper presents a new method to adapt the translation phrase-table of the SMT system. Our experiments were conducted for the English-Vietnamese language pair in the direction from English to Vietnamese. We use specific domain corpus comprise of two specific domains: Legal and General. The data has been collected from documents, dictionaries and the IWSLT 2015 organisers for the English-Vietnamese translation task.

In our works, we train a translation model with parallel corpus in general domain, then we train a probability classifier model with monolingual corpus in legal domain, we use the classification probability of phrase on target side of phrase translation table to recompute the direct translation probability of the phrase translation table. This is the first adaptation method for the phrase translation table of the SMT system, especially for low-resource language pairs as English-Vietnamese language pair. For comparison, we train a baseline SMT system and a Neural Machine Translation system (NMT) to compare with our method. Experimental results showed that our method significantly outperforms the baseline system. Our system improved the translation quality of the machine translation system on the out-of-domain data (*legal domain*) up to 0.9 BLEU points compared to the baseline system. Our method has also been accepted for presentation at the 31st Asia Pacific conference on language, information and computation.

The paper is organized as follows. In the next section, we present related works on the problem of adaptation in SMT; Section 3 describes our method; Section 4 describes and discusses the experimental results. Finally, we end with a conclusion and the future works in Section 5.

2. Related works

Domain adaptation for machine translation is known to be a challenging research problem that has substantial real-world application and this has been one of the topics of increasing

interest for the recent years. Recently, the studies of domain adaptation for machine translation have focused on data-centric or model-centric. Some authors used out-of-domain monolingual data to adapt the language model. The main advantage of language model adaptation in contrast with translation model adaptation, these methods use only monolingual out-of-domain data.

For many language pairs and domains, no new-domain parallel training data is available. In [14] machine translate new-domain source language monolingual corpora and use the synthetic parallel corpus as additional training data by using dictionaries and monolingual source and target language text.

In [5] build several specific domain translation systems, then train a classifier model to assign the input sentence to a specific domain and use the specific domain system to translate the corresponding sentence. They assume that each sentence in test set belongs to one of the already existing domains.

In [11] build the MT system for different domains, it trains, tunes and deploys a single translation system that is capable of producing adapted domain translations and preserving the original generic accuracy at the same time. The approach unifies automatic domain detection and domain model parameterization into one system.

In [15] used a source classification document to classify an input document into a domain. This work makes the translation model shared across different domains.

Above related works automatically detected the domain and the classifier model works as a “switch” between two independent MT decoding runs.

There are many studies of domain adaptation for SMT, data-centric methods usually focus on selecting training data from out-of-domain parallel corpus and ignoring out-of-domain monolingual data, which can be obtained more easily.

Our method has some differences from above methods. For adapting to the translation phrase-table of SMT system, we build a probability classifier model to estimate the

classification probability of phrases on target side of the translation phrase-table. Then we use these classification probabilities to recompute the direct phrase translation probability $\phi(e|f)$.

3. Our method

In phrase-based SMT, the quality of the SMT system depends on training data. SMT systems are usually trained on large amounts of the parallel corpus. Currently, high-quality parallel corpora of sufficient size are only available for a few language pairs. Furthermore, for each language pair, the sizes of the domain-specific corpora and the number of domains available are limited. The English-Vietnamese is low-resource language pair and thus domains data in this pair are limited, for the majority of domains data, only a few or no parallel corpora are available. However, monolingual corpora for the domain are available, which can also be leveraged.

The main idea in this paper is leveraging out-of-domain monolingual corpora in the target language for domain adaptation for MT. In the phrase-table of SMT system, a phrase in the source language may have many translation hypotheses with a different probability. We use out-of-domain monolingual corpora to recompute the scores of translation probability of these phrases which are defined in out-of-domain.

There are many studies of domain adaptation for SMT, which can be mainly divided into two categories: data-centric and model-centric. Data-centric methods focus on either selecting training data from out-of-domain parallel corpora based on a language model or generating parallel data. These methods can be mainly divided into three categories:

- Using monolingual corpora.
- Synthetic parallel corpora generation.
- Using out-of-domain parallel corpora: multi-domain and data selection.

Most of the related works in section 2 use monolingual corpora to adapt language model or to synthesize parallel corpora, or models selection which are trained with different domains. The English-Vietnamese is low-resource parallel corpora, thus we propose a

new method which only uses monolingual corpora to adapt the translation model by recomputing the score of phrases in the phrase-table and to update the phrase's direct translation probability.

In this section, we first give a brief introduction of SMT. Next, we propose a new method for domain adaptation in SMT.

3.1. Overview of phrase-based statistical machine translation

The figure 1 illustrates the process of phrase-based translation. The input is segmented into a number of sequences of consecutive words (so-called phrases). Each word or phrase in English is translated into a word or phrase in Vietnamese, respectively. Then these output words or phrases can be reordered.

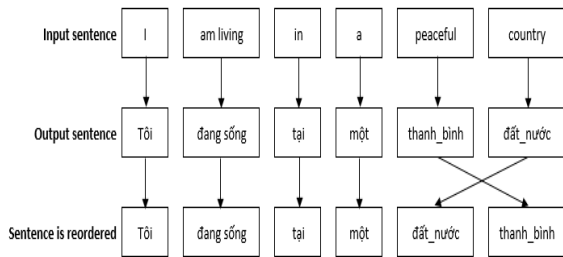


Figure 1. Example illustrates the process of phrase-based translation.

The phrase translation model is based on the noisy channel model [16]. It uses Bayes rule to reformulate the translation probability for translating a input sentence f in one language

into output sentence e in another language. The best translation for a input sentence f is as equation 1:

$$e = \underset{e}{\operatorname{arg\,max}} p(e|f) \quad (1)$$

The above equation consists of two components: A language model assigning a probability $p(e)$ for any target sentence e , and a translation model that assigns a conditional probability $p(e|f)$. The language model is trained with monolingual data in the target language, the translation model is trained with parallel corpus, the parameters of translation model are estimated from a parallel corpus, the best output sentence (e) corresponding to an input sentence (f) is calculated by the after formula 2 and 3.

$$e = \underset{e}{\operatorname{arg\,max}} p(e|f) \quad (2)$$

$$= \underset{e}{\operatorname{arg\,max}} \sum_{m=1}^M \lambda_m h_m(e, f) \quad (3)$$

where h_m is a feature function such as language model, translation model and λ_m corresponds to a feature weight.

The Figure 2 describes the architecture of phrase-based statistical machine translation system. There is some translation knowledge that can be used as language models, translation models, etc. The combination of component models (language model, translation model, word sense disambiguation, reordering model,...).

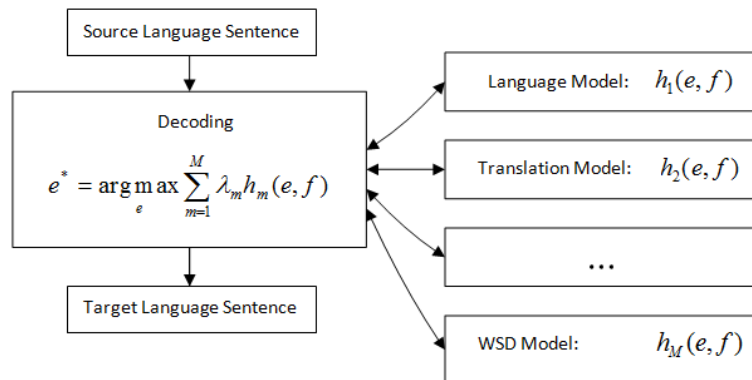


Figure 2. Architecture of phrase-based statistical machine translation.

3.2. Translation model adaptation based on phrase classification

One of the essential parts of our experiments is the classifier used to identify the domain of a target phrase in the phrase-table, the accuracy of the classifier is very important in the final translation score of the sentences from the test set data. The Maximum Entropy was chosen to be the classifier for our experiments.

In this section, we first give an introduction of the maximum entropy classifier. Next, we describe our method for domain adaptation in SMT.

3.2.1. The Maximum Entropy classifier

To build a probability classification model, we use the Stanford classifier toolkit¹ with standard configurations. This toolkit uses a maximum entropy classifier with character n-grams features,... The maximum entropy classifier is a probabilistic classifier which belongs to the class of exponential models. The maximum entropy is based on the principle of maximum entropy and from all the models that fit training data, select the one which has the largest estimate probability. The maximum entropy classifier is usually used to classify text and this model can be shown in the following formula:

$$p(y|x) = \frac{\exp(\sum_k \lambda_k f_k(x, y))}{\sum_k \exp(\sum_k \lambda_k f_k(x, z))} \quad (4)$$

where λ_k are model parameters and f_k are features of the model [17].

We trained the probability classification model with 2 classes which are Legal and General. After training, the classifier model was used to classify a list of phrases in the phrase-table in target side, we consider these phrases to be in the general domain at the beginning. The output of the classification task is a probability of phrase in each domain ($P(\text{legal})$ and $P(\text{general})$), some results of the classification task as in the Figure 3.

phrases in p-table	P(legal)	P(general)
tội_phạm	0.991	0.009
hợp_pháp	0.551	0.449
pháp_lý	0.891	0.109
biểu_tượng	0.519	0.481
bộ_phận	0.688	0.312
cảnh_sát	0.986	0.014
hiệp_hội	0.633	0.367
hậu_quả	0.977	0.023
thẩm_quyền	0.986	0.014
tình_vì	0.870	0.130
tình_huống	0.642	0.358
ảnh_hưởng	0.742	0.258
xử_lý	0.996	0.004
buộc_tội	0.951	0.049
hạn_chế	0.840	0.160
tình_huống	0.642	0.358
bất_hợp_pháp	0.996	0.004
phạm_pháp	0.930	0.070
trái_phép	0.690	0.310
thực_thì	0.976	0.024
thì_hành	0.938	0.062

Figure 3. Some results of the classification task.

3.2.2. Phrase classification for domain adaptation in SMT

The State-of-the-art SMT system uses a log-linear combination of models to decide the best-scoring target sentence given a source sentence. Among these models, the basic ones are a translation model $P(e|f)$ and a target language model $P(e)$.

The translation model is a phrase translation table; this table is a list of the translation probabilities of a specified source phrase f into a specified target phrase e , including phrase translation probabilities in both translation directions, the example about the structure of phrase translation table as the Figure 4.

confirm		khăng_định		0.0571429	0.0238095	0.769	0.142857
confirm		xác_minh_được_không		1	0.0370849	0.2	0.000728863
confirm		xác_nhận		0.0625	0.09375	0.2	0.214286
confirm		xác_nhận_được		1	0.0469315	0.2	0.0306122
consequences		hậu_quả		0.240506	0.259494	0.965	0.465909
consequences		hệ_quả		0.151515	0.206897	0.106383	0.136364
consequences		kết_quả		0.00448431	0.0088889	0.0425532	0.0909091
copyright		bản_quyền		0.543478	0.532258	0.994	0.66
copyright		bảo_hộ		0.142857	0.037037	0.027027	0.02
crime		phạm_tội		0.454545	0.181818	0.949	0.0930233
crime		tội_lỗi		0.0238095	0.0178571	0.0149254	0.0116279
crime		tội_phạm		0.381579	0.275862	0.994	0.372093
crime		tội_ác		0.40625	0.348837	0.19403	0.174419

Figure 4. Example of phrase translation scores in phrase-table.

¹ <https://nlp.stanford.edu/software/classifier.html>

In the Figure 4, the phrase translation probability distributions $\phi(f|e)$ and $\phi(e|f)$, lexical weighting for both directions. Currently, four different phrase translation scores are computed:

1. Inverse phrase translation probability $\phi(f|e)$.
2. Inverse lexical weighting $\text{lex}(f|e)$.
3. Direct phrase translation probability $\phi(e|f)$.
4. Direct lexical weighting $\text{lex}(e|f)$.

In this paper, we only conduct the experiments with translation direction from

English to Vietnamese, thus we only investigate the direct phrase translation probability $\phi(e|f)$ of the phrase-table, the translation hypothesis is higher probability $\phi(e|f)$ value, that translation hypothesis is often chosen more than another, so we use the probability classification model to determine the classification probability of a phrase in the phrase-table, then we recompute the translation probability of phrase $\phi(e|f)$ of this hypothesis based on the classification probability.

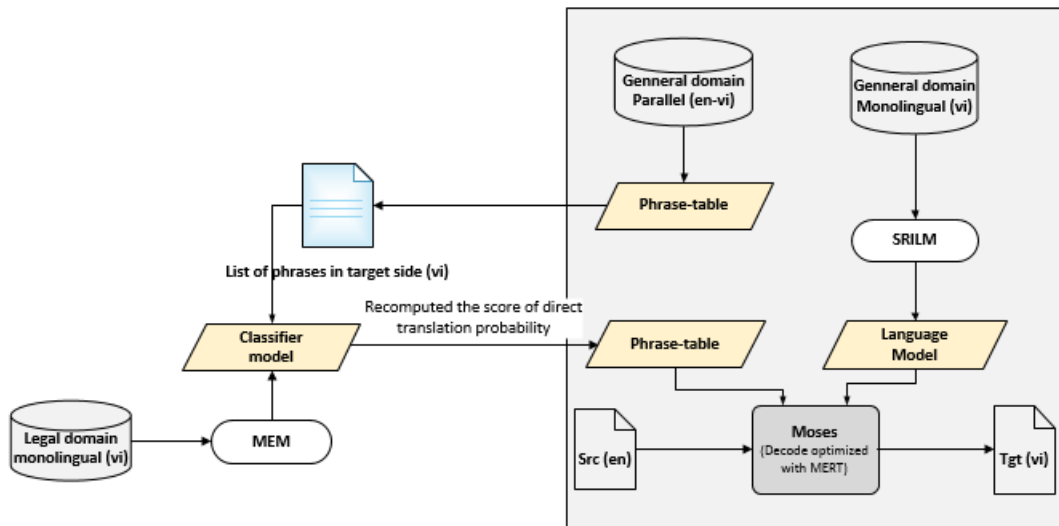


Figure 5. Architecture of the our translation model adaptation system.

Our method can be illustrated in the Figure 5 and summarized by the following:

1. Build a probability classification model (using the maximum entropy classifier with two classes, legal and general) with monolingual data on legal domain in Vietnamese.
2. Training a baseline SMT system with parallel corpus on general domain with translation direction from English to Vietnamese.
3. Extract phrases on target side of the phrase-table of the baseline SMT system and using the probability classification model for these phrases.

4. Recompute the direct translation probability $\phi(e|f)$ of phrases of the phrase-table for phrases are classified into the legal label.

4. Experimental Setup

In this section, we describe experimental settings and report empirical results.

4.1. Data sets

We conduct experiments on the data sets of the English-Vietnamese language pair. We consider two different domains that are legal domain and general domain. Detailed statistics for the data sets are given in the Table 1.

Out-of-domain data: We use monolingual data on legal domain in the Vietnamese language, this data set is collected from documents, dictionaries,... consists of 2238 phrases, manually labelled, including 526 in-of-domain phrases (*in legal domain and label is lb_legal*) and 1712 out-of-domain phrases (*in general domain and label is lb-general*). Here

the phrase concept is similar to the phrase concept in the phrase translation table, this concept means nothing more than an arbitrary sequence of words, with no sophisticated linguistic motivation. This data set is used to train the probability classification model by the maximum entropy classifier with 2 classes, legal and general.

Table 1. The Summary statistical of data sets: English-Vietnamese

Data Sets		Language	
		English	Vietnamese
Training	Sentences	122132	
	Average Length	15.93	15.58
	Words	1946397	1903504
	Vocabulary	40568	28414
Dev	Sentences	745	
	Average Length	16.61	15.97
	Words	12397	11921
	Vocabulary	2230	1986
General-test	Sentences	1046	
	Average Length	16.25	15.97
	Words	17023	16889
Legal-test	Vocabulary	2701	2759
	Sentences	500	
	Average Length	15.21	15.48
	Words	7605	7740
	Vocabulary	1530	1429

Additionally, we use 500 parallel sentences on legal domain in English-Vietnamese pair for test set.

In-of-domain data: We use the parallel corpora sets on general domain to training SMT system. These data sets are provided by the IWSLT 2015 organisers for the English-Vietnamese translation task, consists of 122132 parallel sentences for the training set, 745 parallel sentences for development set and 1046 parallel sentences for the test set.

Preprocessing: Data preprocessing is very important in any data-driven method. We carried out preprocessing in two steps:

• **Cleaning Data:** We performed cleaning in two phases, *phase-1*: Following the cleaning

process described in [18] and *phase-2*: Using the corpus cleaning scripts in Moses toolkit [19] with minimum and maximum number of tokens set to 1 and 80 respectively.

• **Word Segmentation:** In English, whitespaces are used to separate words [20] but Vietnamese does not have morphology [21] and [20]. In Vietnamese, whitespaces are not used to separate words. The smallest meaningful part of Vietnamese orthography is a syllable [22]. Some examples of Vietnamese words are shown as follows: *single words* "nhà" - house, "nhặt" - pick up, "mua" - buy and "bán" - sell. *Compound words*: "mua-bán" - buy and sell, "bàn-ghế" - table and chair, "cây-cối" - trees, "đường-xá" - street, "hành-chính" -

administration. Thus, a word in Vietnamese may consist of several syllables separated by whitespaces.

We used vntokenizer toolkit [23] to segment for Vietnamese data sets, this segmentation toolkit is quite popular for Vietnamese and we used tokenizer script in Moses to segment for English data sets.

4.2. Experiments

We performed experiments on the Baseline-SMT and Adaptaion-SMT systems:

- The Baseline-SMT is a SMT baseline system. This system is the phrase-based statistical machine translation with standard settings in the Moses toolkit² [24], this is a state-of-the-art open-source phrase-based SMT system. In our systems, the weights of feature functions were optimized using MERT [25]. The Baseline-SMT is trained on the general domain (*in-of-domain*) data set and the Baseline-SMT system is evaluated sequentially on the General-test and Legal-test data sets.

Source sentences (on the Legal_domain)	Target sentences				Reference sentences
	Baseline_SMT system	Adaptation_SMT system	NMT system	Google Translate	
the working party took note of this commitment .	bữa tiệc làm_việc nốt_nhạc đã cam_kết này .	nhóm làm_việc ghi_nhận cam_kết này .	buổi tiệc đã nhận được sự cam_kết về những kiểm_khuyết của sự cam_kết này .	nhóm làm việc đã lưu_ý về cam_kết này.	ban công tác đã ghi_nhận cam_kết này .
according to the general statistical office , services had accounted for 37.98 percent of vietnam 's gdp in 2004	theo tổng_quát văn_phòng thông_kê , dịch_vụ đã chiếm 37.98% của việt_nam là gdp vào năm 2004 .	theo tổng_cục thông_kê , dịch_vụ đã chiếm 37.98% của việt_nam là gdp vào năm 2004 .	theo các văn_phòng thông_kê , dịch_vụ đã có những phần_trăm trong số gdp của việt nam trong sự kim_kệp của gdp ở dorset .	Theo cơ quan thông_kê chung , các dịch vụ đã chiếm 37,98% gdp của Việt Nam năm 2004.	theo tổng_cục thông_kê , dịch_vụ chiếm 37,98% gdp năm 2004 của việt_nam .
renewable certificates valid for five years were granted by the construction departments of cities and provinces .	giấy chứng nhận tái_tạo có giá_trị trong 5 năm qua là cấu_trúc của các thành_phố và đại_lục .	giấy_phép gia_hạn có giá_trị trong 5 năm bởi cấu_trúc của các thành_phố và đại_lục .	những mẫu giấy tái_tạo có giá_trị trong 5 năm là do các nhà hoạt_động_xây_dựng của các thành_phố và các quận biên_động .	Giấy chứng nhận tái_tạo có giá trị trong năm năm được cấp bởi các sở xây dựng của thành phố và các tỉnh.	sở xây dựng các tỉnh và thành_phố cấp giấy_phép hành_nghề có hiệu_lực 5 năm và các giấy_phép này có_thể được gia_hạn .
the economic police received specialized training on intellectual property enforcement .	cảnh_sát kinh_tế được đào_tạo chuyên về cơ_quan sở_hữu_trí_tuệ .	cảnh_sát kinh_tế được đào_tạo chuyên về thực_thi quyền_sở_hữu_trí_tuệ .	cảnh_sát được đào_tạo chuyên về các ngăn vi_phạm sở_hữu_trí_tuệ .	cảnh sát kinh tế được đào tạo chuyên môn về thực thi sở hữu trí tuệ.	cảnh_sát kinh_tế được đào_tạo chuyên_sâu về thực_thi quyền_sở_hữu_trí_tuệ .
administrative measures only applied to acts of low gravity .	đo_hành_chính chỉ áp_dụng cho hành_động của trọng_lực thấp .	các biện_pháp hành_chính chỉ áp_dụng cho các hành_vi nhghiêm_trọng thấp .	các biện_pháp quản_lý chỉ áp_dụng vào những hành_động thấp của những thiên_thể thấp .	biện pháp hành chính chỉ áp dụng cho các hành vi trọng lực thấp.	các biện_pháp hành_chính chỉ áp_dụng với những hành_vi có tính nhghiêm_trọng thấp .
evidence collected during an administrative procedure could be used by the civil court if necessary in accordance with civil procedure code of 2004 .	bằng_chứng thu_thập được trong một ca_hành_chính có_thể được sử_dụng bởi những tòa dân_sự nếu cần_thiết theo thủ_tục dân_sự của năm 2004 .	bằng_chứng thu_thập được trong thủ_tục hành_chính có_thể được sử_dụng bởi tòa dân_sự nếu cần_thiết theo bộ_luật thủ_tục dân_sự của năm 2004 .	bằng_chứng được thu_thập trong một thủ_tục quản_lý có_thể được sử_dụng bởi tòa_án dân_sự nếu cần_thiết trong hệ_thống dân_sự của năm 2004 .	bằng chứng thu thập trong một thủ_tục hành chính có thể được sử dụng bởi tòa án dân sự nếu cần thiết theo quy tắc tố tụng dân sự năm 2004.	chứng_cứ thu được trong quá trình xử lý hành_chính sẽ được sử_dụng tại tòa dân_sự nếu thấy cần_thiết theo bộ_luật tố_tụng dân_sự năm 2004 .

Table 3. Some examples in our experiments.

- The Adaptation-SMT is based on the Baseline-SMT system after being adapted to the translation model by recomputing the direct translation probability $\phi(e|f)$ of phrases in the phrase translation table, the Adaptaion-SMT is evaluated on the Legal-test data set².

We train a language model with 4-gram and Kneser-Ney smoothing was used in all the experiments. We used SRILM³ [26] as the

language model toolkit. For evaluate translation quality of the Baseline-SMT system and Adaptaion-SMT system, we use the BLEU score [27].

For comparison, we also built a Neural Machine Translation (NMT) system use OpenNMT toolkit⁴ [28], the NMT system is trained with the default model, which consists of a 2-layer LSTM with 500 hidden units on both the encoder/decoder.

² <http://www.statmt.org/moses/>

³ <http://www.speech.sri.com/projects/srilm/>

⁴ <http://opennmt.net/>

4.2.1. Results

Table 2. The experiment results of the Baseline-SMT system and Adaptaion-SMT system

SYSTEM	BLEU SCORE
Baseline-SMT (<i>General-test</i>)	31.3
Baseline-SMT (<i>on Legal-test</i>)	28.8
Adaptaion-SMT (<i>on Legal-test</i>)	29.7
Baseline-NMT (<i>on General-test</i>)	30.1
Baseline_NMT (<i>on Legal-test</i>)	20.9

The Table 2 showed that the baseline systems (*the SMT and NMT system*) are trained on the general domain data set, if the test data set (*here is the General-test data set*) is in the same domain as the training data, the BLEU score will be 31.3 for the Baseline-SMT system and 30.1 for the Baseline-NMT system. If the test data set is on the legal domain (*here is the Legal-test data set*), the BLEU score will be

28.8 for the Baseline-SMT system and 20.9 for Baseline-NMT system.

The Table 2 also showed that the SMT system is trained on the general domain if the test domain is different from the training domain, the quality of the translation quality is reduced. In these experiments, the BLEU score was reduced by 2.5 points from 31.3 to 28.8. The Adaptaion-SMT system is adapted by our technique will improve the quality of the translation system. In these experiments, the BLEU score is improved to 0.9 points from 28.8 up to 29.7.

The experiment results also showed that the SMT system has better results than the NMT system when translation systems are trained with the same low-resource domains of English-Vietnamese language pair such as legal domain and some other domains.

4.2.2. Analysis and discussion

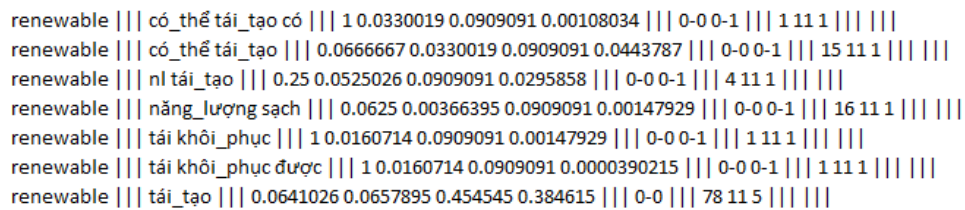


Figure 6. Examples about the direct translation probability of this phrase in phrase-table.

Some examples in the Table 3, when systems translate source sentences in legal domain from english to vietnamese language. In the third sentence, the phrase “renewable” in context “renewable certificates valid for five years were granted by the construction departments of cities and provinces” (*source sentence column*) should be translated into “gia-hạn” as reference sentence but the Baseline-SMT system has translated the phrase “renewable” into “tái-tạo”, the NMT system has translated that phrase into “tái-tạo”, the Google Translate has translated that phrase into “tái tạo” and the Adaptaion-SMT system has translated the phrase “renewable” into “gia-hạn” like reference sentence.

The first, the Baseline-SMT system has translated the phrase “renewable” into “tái-

tạo” because the direct translation probability (*4th column in Figure 6*) of this phrase in phrase-table of Baseline-SMT system is highest (0.454545), and the direct translation probability into “gia-hạn” is lower (0.0909091). Therefore, when the SMT system combines component models as formulas 1, the ability to translate into “tái-tạo” will be higher “gia-hạn”.

Later, apply the phrase classification model to compute the probability of “gia-hạn” and “renewable” phrase in legal domain, the probability of “gia-hạn” is higher than that, then update this value to phrase-table and the direct translation probabilities $\phi(e|f)$ of phrase are recomputed. Therefore, the Adaptation-SMT has translated “renewable” phrase into “gia-hạn”

Some other examples in the Table 4.2 showed that translation quality of Adaptaion-SMT system is better than the Baseline-SMT system and with low-resource translation domains in English-Vietnamese language, the SMT system has more advantages than the NMT system.

5 Conclusions and future works

In this paper, we presented a new method for domain adaptation in Statistical Machine Translation for low-resource domains in English-Vietnamese language pair. Our method only uses monolingual out-of-domain data to adapt the phrase-table by recomputing the phrase's direct translation probability $\phi(e|f)$. Our system obtained an improved on the quality of machine translation in the legal domain up to 0.9 BLEU points over baseline. Experimental results show that our method is effective in improving the accuracy of the translation.

In future, we intend to study this problem with other domains, the benefits of word embedding in phrase classification and integrate automatically our technique to decode of SMT system.

References

- [1] Philipp Koehn, Franz Josef Och, Daniel Marcu, Statistical phrase-based translation, In Proceedings of HLT-NAACL, Edmonton, Canada, 2003, pp. 127-133.
- [2] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes and Jeffrey Dean, Google's neural machine translation system: Bridging the gap between human and machine translation, CoRR, abs/1609.08144, 2016.
- [3] Luisa Bentivogli, Arianna Bisazza, Mauro Cettolo and Marcello Federico, Neural versus phrase-based machine translation quality: A case study, 2016.
- [4] Barry Haddow, Philipp Koehn, Analysing the effect of out-of-domain data on smt systems, In Proceedings of the Seventh Workshop on Statistical Machine Translation, 2012, pp. 422-432.
- [5] Boxing Chen, Roland Kuhn and George Foster, Vector space model for adaptation in statistical machine translation, Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, 2013, pp. 1285-1293.
- [6] Daniel Dahlmeier, Hwee Tou Ng, Siew Mei Wu, Building a large annotated corpus of learner english: The nus corpus of learner english, In Proceedings of the NAACL Workshop on Innovative Use of NLP for Building Educational Applications, 2013.
- [7] Eva Hasler, Phil Blunsom, Philipp Koehn and Barry Haddow, Dynamic topic adaptation for phrase-based mt, In Proceedings of the 14th Conference of the European Chapter of The Association for Computational Linguistics, 2014, pp. 328-337.
- [8] George Foster, Roland Kuhn, Mixture-model adaptation for smt, Proceedings of the Second Workshop on Statistical Machine Translation, Prague, Association for Computational Linguistics, 2007, pp. 128-135.
- [9] George Foster, Boxing Chen, Roland Kuhn, Simulating discriminative training for linear mixture adaptation in statistical machine translation, Proceedings of the MT Summit, 2013.
- [10] Hoang Cuong, Khalil Sima'an, and Ivan Titov, Adapting to all domains at once: Rewarding domain invariance in smt, Proceedings of the Transactions of the Association for Computational Linguistics (TAACL), 2016.
- [11] Ryo Masumura, Taichi Asam, Takanobu Oba, Hirokazu Masataki, Sumitaka Sakauchi, and Akinori Ito, Hierarchical latent words language models for robust modeling to out-of domain tasks, Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015, pp. 1896-1901
- [12] Chenhui Chu, Raj Dabre, and Sadao Kurohashi. An empirical comparison of simple domain adaptation methods for neural machine translation, 2017.
- [13] Markus Freitag, Yaser Al-Onaizan, Fast domain adaptation for neural machine translation, 2016.
- [14] Jia Xu, Yonggang Deng, Yuqing Gao and Hermann Ney, Domain dependent statistical

- machine translation, In Proceedings of the MT Summit XI, 2007, pp. 515-520.
- [15] Hua Wu, Haifeng Wang Chengqing Zong, Domain adaptation for statistical machine translation with domain dictionary and monolingual corpora, In Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008), Manchester, UK, 2008, pp. 993-1000.
- [16] Adam Berger, Stephen Della Pietra, and Vincent Della Pietra, A maximum entropy approach to natural language processing, *Computational Linguistics*, 22, 1996.
- [17] Santanu Pal, Sudip Naskar, Josef Van Genabith, Uds-sant, English-German hybrid machine translation system, In Proceedings of the Tenth Workshop on Statistical Machine Translation, Lisbon, Portugal, September, Association for Computational Linguistics, 2015, pp. 152-157.
- [18] Louis Onrust, Antal van den Bosch, Hugo Van hamme, Improving cross-domain n-gram language modelling with skipgrams, Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, Berlin, Germany, 2016, pp. 137-142.
- [19] Mark Aronoff, Kirsten Fudeman, *What is morphology*, Vol. 8. John Wiley and Sons, 2011.
- [20] Laurence C. Thompson, The problem of the word in vietnamese, In *Journal of the International Linguistic Association* 19(1) (1963) 39-52. <https://doi.org/10.1080/00437956.1963.11659787>.
- [21] Binh N. Ngo, The Vietnamese language learning framework, *Journal of Southeast Asian Language Teaching* 10 (2001) 1-24.
- [22] Le Hong Phuong, Nguyen Thi Minh Huyen, Azim Roussanaly, Ho Tuong Vinh, A hybrid approach to word segmentation of vietnamese texts, 2008.
- [23] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, Evan Herbst, Moses: Open source toolkit for statistical machine translation, In *ACL-2007: Proceedings of demo and poster sessions*, Prague, Czech Republic, 2007, pp.177-180.
- [24] Franz Josef Och, Minimum error rate training in statistical machine translation, In *Proceedings of ACL*, 2003, pp.160-167.
- [25] Andreas Stolcke, Srilm - an extensible language modeling toolkit, in *proceedings of international conference on spoken language processing*, 2002.
- [26] Papineni, Kishore, Salim Roukos, Todd Ward, WeiJing Zhu, Bleu: A method for automatic evaluation of machine translation, *ACL*, 2002.
- [27] G. Klein, Y. Kim, Y. Deng, J. Senellart, A.M. Rush, OpenNMT: Open-Source Toolkit for Neural Machine Translation. ArXiv e-prints.
- [28] Pratyush Banerjee, Jinhua Du, Baoli Li, Sudip Kr. Naskar, Andy Way and Josef van Genabith, Combining multi-domain statistical machine translation models using automatic classifiers, In *Proceedings of AMTA 2010.*, 2010.