

Building a Statistical Machine Translation System for Translating Patent Documents

Jeff Ma

50 Moulton St
Cambridge, MA 02138
USA

jma@bbn.com

Spyros Matsoukas

50 Moulton St
Cambridge, MA 02138
USA

smatsouk@bbn.com

Abstract

This paper describes the work we conducted for building a statistical machine translation (SMT) system for the Chinese-English sub-task of the NTCIR-9 patent MT evaluation. Our results show that most of the generic techniques we had developed for improving SMT performance work on patent data as well, and the changes we made to our SMT system training procedure in order to address special characteristics of patent documents produced additional improvements.

1 Introduction

Compared to human translation, machine translation (MT) is fast and low-cost, so it has been used for translating large numbers of patent documents. Patent documents are juridical documents, which are typically more structured than general documents, and they have their own special characteristics. People tried to utilize these special characteristics in various applications, such as categorization of patent documents in (Kim and Choi, 2007; Iwayama et al., 2003), and machine translation of patent documents in (Shimohata, 2005; Ofersgaard and Povlsen, 2007; Jin, 2010). Some of the patent document characteristics make MT easier, for example, the presence of well-structured sentences and less ambiguity of word meanings. On the other hand, some characteristics become challenges for MT, for example, long and complicated sentence structures, technical terminology

and new terms that are originally defined by patent applicants. Due to these challenges people have explored various strategies for improving patent MT quality, such as combining SMT with rule-based MT in (Terumasa, 2007; Wang, 2009; Jin, 2010), with promising results. In our work, we mainly focused on building an SMT system for translating Chinese patent documents to English and designing techniques for the SMT system to handle special characteristics of patent data better.

The paper is organized as follows: Section 2 briefly describes our SMT training procedure; Section 3 presents data preparation for building the patent MT system; and Section 4 reports incremental gains from a few methods we implemented specially for patent MT.

2 Our SMT system

We have been building our SMT systems based on the model described in (Shen et al., 2008), which employs hierarchical rules to translate strings in the source language to dependency trees in the target language. Recently we improved the SMT model with two techniques: use of a large number (50,000) of features, similar to the method reported in (Chiang et al., 2009), and discriminative training of feature weights to maximize the expected BLEU (Rosti et al., 2010). Since the expected BLEU criterion is continuous and differentiable, gradient descent may be performed, thus supporting weight tuning for a large number of features. The use of the 50,000 features yielded gains similar to those reported in (Chiang et al., 2009). For convenience,

we will refer to the features we used prior to adding the 50,000 features as the ‘regular’ features. We use GIZA++ (Och and Ney, 2003) for training word alignment models.

We have built various Chinese-English SMT systems. One of them is a “newswire” MT system for translating Chinese newswire text. We used this system to help set up a development set and a test set for building the patent SMT system. The “newswire” MT system was trained on a parallel training corpus that includes 227 million (227M) words, the majority of which is newswire text. The collections in this corpus had been released by the Linguistic Data Consortium (LDC) for the DARPA GALE project. Our MT system uses a tri-gram target language model (LM) to generate n-best hypotheses and then ranks the n-best with a 4-gram LM. We build n-gram LMs with the modified Kneser-Ney smoothing (Chen and Goodman, 1998). For the “newswire” MT system we trained the target (English) LMs on an English corpus that consists of more than 6 billion (6B) words of news text, out of which 227M words come from the English side of the Chinese-English parallel corpus, 2.2 billion words from the 4th edition of the LDC English Gigaword monolingual data release, 2.5 billion words from Google news and 1.6 billion words from news text that we downloaded from various websites, such as BBC, Xinhua News, and The Arab News. We denote this LM as “6B-nw-LM”. This “newswire” MT system produced a BLEU (Papineni et al., 2002) score of 26.22 on the GALE Phase 4 Chinese newswire evaluation test set that includes 490 sentences – with one reference translation per sentence.

3 Patent data preparation

For building the patent MT system, we used the data released by the NTCIR-9¹ evaluation organizers for the Chinese-English sub-task of the NTCIR-9 patent MT evaluation. The released data includes a parallel training corpus that consists of one million (1M) Chinese-English sentence pairs and a development data set that consists of two thousand (2K) bilingual sentence pairs.

The 1M parallel sentence pairs were extracted from patents that were published before 2006. The total number of English words in the 1M parallel

sentence pairs was close to 45 million (45M). The 2K sentence pairs in the development set were extracted from 103 patents that were published in the years of 2006 and 2007. The full contents of the 103 patents, including titles, abstracts and full descriptions, were also provided in both Chinese and English languages. As the evaluation organizers stated, people could use the full contents – referred to as “context data” – for building MT systems. We explored using this context data to adapt the target LM. We split this development set into two subsets, one for tuning the MT system and one for measuring the performance. To make the two subsets similar in terms of translation difficulty, we first translated these 2K sentences with the “newswire” MT system, and then split the 103 patent documents into two subsets, roughly half-and-half, based on their translation error rate (TER) (Snover et al., 2010), resulting in two subsets of approximately equal TER scores. With this splitting, we ended up with 1039 sentences from 54 patents in the tuning set – denoted as “Tune” in this paper – and 961 sentences from 49 patents in the test set – denoted as “Test”. The mixed-case BLEU scores measured on the whole 2K development set, the Tune, and the Test sets are listed in Table 1. The scores on the Tune and Test are close.

Data set	2K-dev	Tune	Test
BLUE	15.38	15.87	14.77

Table 1. BLEU scores measured on the 2K development set, the Tune, and the Test sets using the “newswire” MT system

Recall that the “newswire” system produced a BLEU score of 26.22 on the newswire evaluation test set. However, on the 2K patent sentences it performed significantly worse – a BLEU score of 15.38. We re-tuned the decoding parameters for the “newswire” MT system with the new Tune set and the re-tuning only improved the BLEU on the Test set slightly –from 14.77 to 15.64. This implies that this big performance degradation mainly resulted from mismatches between the training and test data.

Besides the parallel training corpus and the development data set, the NTCIR-9 committee also released a monolingual English patent corpus for the purpose of training English LMs. This corpus includes US patent documents published in the

¹ <http://ntcir.nii.ac.jp/PatentMT/>

period 1993-2005, totaling 14 billion (14B) words. We used this corpus for training our English LMs.

Since we focused on the mixed-case performance in our work, we will report only the mixed-case BLEU scores measured on the Test set for all experiments shown below, unless specified otherwise.

4 Building a patent SMT system

4.1 Training the MT system

We first re-trained the MT model with the 45M word patent parallel corpus. Before training the word alignment models, we segmented words in the Chinese sentences with a 52K lexicon by using a left-to-right and longest-match-first algorithm, which generated 41 million (41M) Chinese words in the 1M sentences.

We trained two sets of English LMs. One was trained with only the 45M English words from the 1M parallel corpus and the other one with the 45M words plus the 14B monolingual English corpus. We denote the former one as “45M-pt” LM and the latter one as “14B-pt” LM.

We looked into effects of the three different LMs, “6B-nw”, “45M-pt” and “14B-pt”, on the MT performance. Mixed-case BLEU scores of the re-trained MT model with the three LMs are listed in Table 2.

MT model	227M newswire	45M patent	45M patent	45M patent
LM	6B-nw	6B-nw	45M-pt	14B-pt
Test	14.77	30.71	34.01	36.16

Table 2. Effect (BLEU scores) of the three LMs when used with the patent SMT system

As can be seen, the use of the 14B monolingual patent data in the LM training helped improve the performance by about 2 BLEU points (from 34.01 to 36.16).

In the above set of experiments we used only the regular features (not including the 50K features). The main reason was to save time for exploring the best strategies to build the patent MT system. For the same reason we also used the smaller LM – “45M-pt” – in many of our following investigation experiments. Unless specified otherwise, experiments reported later in this paper used the “45M-pt” LM and the regular features.

Hence, the system trained with the 45M parallel patent data (the 4th column in Table 2) serves as a baseline for our later efforts to improve the patent MT performance.

4.2 Addressing special characteristics of patent data

We first saw that, compared to the Chinese news-wire text data, the Chinese patent text includes significantly more special strings that are not written in Chinese characters, such as English words, patent numbers, mathematical expressions and abbreviation names for materials. This is one characteristic of the patent data. Since all the special strings are written in ASCII characters, we call them ASCII strings. We found many of the ASCII strings occurring in the 45M word patent parallel corpus were not aligned properly during the word alignment training. The main reason was inconsistent tokenization of the ASCII strings on the source and target sides. For example, the ASCII string “IS-1000” was tokenized as itself when occurring in the Chinese sentences but tokenized as “IS – 1000” when occurring in the English sentences. To remove such inconsistency we tokenized the ASCII strings in the Chinese sentences in the same way as we tokenized the English sentences. The system trained with this consistent tokenization of ASCII strings is denoted as “+ consistent tokenization” in Table 3, where we use the sign “+” to indicate changes applied on top of the system shown in the preceding row. Compared to the “baseline”, the consistent tokenization of the ASCII strings improved the performance by about half a BLEU point.

The second thing we did was to increase sharing of translation rules and LM n-gram scores among certain types of special tokens. When training Chinese-English MT systems, we let “infrequent” numbers – all numbers except numbers in the range 1-31 – share translation rules and LM n-gram scores. The sharing mechanism is as follows:

- 1) Train word alignment models after replacing “infrequent” numbers on both sides of the parallel corpus with a special “number token”
- 2) Train LMs after applying the same number replacement on the LM training corpus

- 3) Before translating test sentences, conduct the same number replacement on test sentences and save replacement information that includes the original numbers and their places in the sentences
- 4) After translating the test sentences, replace the special number tokens occurring in the MT hypotheses with their corresponding original numbers based on the number replacement information and source-to-target word alignment information that the MT decoder outputs during the translation

This sharing mechanism improves our MT performance. Because there are more special tokens in the patent data, such as patent identification numbers and mathematical expressions, we applied the translation rule and n-gram sharing mechanism on 4 more special tokens:

1. patent identification numbers – all 7-digit whole numbers, such as 5,716,812 and 5869649
2. name abbreviations – ASCII strings occurring in the Chinese sentences that consist of only English characters and digits, such as “PMMA” and “CO2”
3. numbers with labels – numbers followed with the commonly-used unit labels, such as “1.03ml” and “20.8g”
4. math expressions – items that consists of any of the math signs, such as “x=0.25” and “a+b”

We applied this special token rule and n-gram score sharing on top of the “+ *consistent tokenization*” system, this new system is denoted as “+ *more sharing*” in Table 3. As can be seen, the rule and LM n-gram sharing on the 4 special tokens produced a 0.4 gain on the BLEU.

System	Test
<i>Baseline</i>	34.01
+ <i>consistent tokenization</i>	34.56
+ <i>more sharing</i>	34.97
+ <i>patent case-LM</i>	36.47
+ <i>optimized word segmentor</i>	36.95

Table 3. Improvements (on BLEU) from addressing patent data related issues

4.3 Re-training the casing LM

We case our MT outputs with a tri-gram LM that is trained with mix-cased English text. The casing algorithm searches among all the possible casing combination of the words in a sentence for the path that has the highest likelihood against the tri-gram casing LM. Our initial casing LM was trained on the mixed-case version of the newswire 6B LM training corpus. As shown before, the newswire data differs significantly from the patent data in terms of the data characteristics. Therefore, we re-trained the casing LM with the mixed-case English sentences from the 45M patent parallel corpus. This new casing LM improved the mixed-case BLEU score by 1.5 points, as shown in the row “+ *patent case-LM*” in Table 3.

4.4 Optimizing the Chinese word segmentor

In the experiments we have reported so far, we segmented the Chinese words with a 52K Chinese word lexicon by using a simple left-to-right and longest-match-first algorithm. The 52K lexicon is an optimized subset of a big Chinese word lexicon that includes 121K entries². Our lexicon optimization procedure starts with a big lexicon and gradually removes words from the lexicon that are not aligned well – by measuring if the removal improves the MT performance. The procedure is as follows:

1. Segment Chinese words in the parallel corpus with an initial big word lexicon
2. Train word alignments and measure the MT performance on a test set
3. Remove from the lexicon any words that are aligned less than “*Threshold*” times
4. Segment Chinese words in the parallel corpus with the reduced lexicon
5. Train a new word alignment model
6. Measure MT performance with the new word alignment model
7. If the performance gets improved, go to Step 3 with an increased value for the “*Threshold*”. Otherwise, stop.

On the 227M Chinese-English parallel corpus we started with the 121K word lexicon and ran a few iterations of the optimization by increasing the

² It consists of the words from the Chinese word lexicon released by LDC (LDC96L15) and words we acquired from a few websites.

“*Threshold*” value gradually – from 5 to 10 to 20 and to 30. We obtained the best MT performance when the “*Threshold*” was set to 20 and the lexicon size was reduced to 52K.

We ran the same lexicon optimization on the 45M patent parallel corpus, but starting with a 62K lexicon that includes the 52K lexicon and 10K new words we extracted from the ADSO word translation lexicon³. By increasing the threshold from 2 to 3 to 4 – much smaller values due to the significantly less amount of training data – we got the best MT performance when the lexicon was reduced to 32K (at the threshold = 3). The BLEU score on the Test set with the initial 62K lexicon was 36.07% and was increased to 36.95% with the optimized 32K lexicon. The performance of the system that used this optimized 32K lexicon to segment the Chinese words is shown in the row “+ *optimized word segmentor*” of Table 3. As can be seen, this lexicon optimization improved the MT performance by 0.5 BLEU points, compared to the system that used the 52K lexicon.

4.5 Using more features

We then added more features to the system. As mentioned before, the total number of new features we extracted was 50,000 (50K). These features came from 8 feature categories:

1. Does the rule contain the target phrase X?
2. Does the rule translate word X to word Y?
3. Does the rule translate POS X to POS Y?
4. Was this rule seen exactly once in the training?
5. Do the two non-terminals in source switch position in the target?
6. Does the source word X align to exactly two target words?
7. How often was the lexical source-target pair (X, Y) seen in the training corpus Z?
8. Is the target non-terminal X filled by the target non-terminal Y?

Over-fitting is a well-known problem for tuning weights for a large number of features. We discriminatively trained feature weights to maximize

³ The latest release (v5.077) of the ADSO dictionary consists of 185K entries, which is free to the public (<http://www.adsotrans.com/downloads>). The dictionary includes many phrasal translations. We extracted entries that have only a single word on the English side and treated the tokens on the Chinese side as Chinese words, and then we selected 10K words that are not in the 121K Chinese lexicon.

the expected BLEU by using the same technique as reported in (Rosti et al., 2010). The expected BLEU was computed with the same formula as the BLEU computation in (Papineni, et al., 2002), but the n-gram counts and matches are expected versions that are derived from n-best hypotheses.

System	Tune	Test
<i>optimized word segmentor</i>	38.66	36.95
<i>add 50K features</i>	44.57	37.38
<i>add top-100 best features</i>	42.82	37.71

Table 4. System performance (BLEU scores) with different numbers of new features added

After adding the 50K features, we noticed that the gap between the BLEU scores on the Tune and Test sets got significantly large. As shown in the row “*add 50K features*” in Table 4, the gap between the Tune and Test sets was 7 BLEU points, which is much larger than that we observed when we were conducting the same tuning for the system trained on the 227M parallel corpus. So the over-fitting problem got worse in the tuning here because of the small size of the tuning set. To alleviate the over-fitting we tried to reduce the number of the new features. Based on the tuned weights for the 50K features, we selected the top-100 ones that had the highest weights and then added only these 100 features to the system. This experiment is shown in the row “*add top-100 best features*” of Table 4. As shown, the use of the top-100 best features alleviated the over-fitting and the performance on the Test set improved. Compared to the baseline – the “*optimized word segmentor*” system, the use of the 100 extra features helped produce 0.8 BLEU gain.

4.6 LM adaptation

For MT we adopted an LM adaptation approach, similar to (Snover et al., 2008), that interpolates a general LM with an LM estimated from text data closely related to the test document that is being translated. We acquire the related text data through the cross-lingual information retrieval (CLIR) technique. This LM adaptation has helped improve performance for most of our MT systems. We use the term “bias LM” to refer to the LM estimated from the CLIR-retrieved text. While translating a test document, we compute log LM scores according to,

$$\begin{aligned} & \log(\text{score}_{LM}) \\ &= \log(\text{score}_{\text{generalLM}} + \alpha * \text{score}_{\text{biasLM}}) \end{aligned} \quad (\text{Eqn.1})$$

where “*generalLM*” denotes the general LM and “*biasLM*” the bias LM. “ α ” is an interpolation weight that is document-dependent and automatically estimated. For a test document – d , the adaptation procedure is as follows:

- 1) treat the document as a query and run a CLIR tool to extract N related passages from a large monolingual text corpus in the target language
- 2) compute the mean, $\mu(d)$, and standard deviation, $\sigma(d)$, of the CLIR scores of all the N passages
- 3) select passages whose CLIR scores are higher than $\mu(d) + T * \sigma(d)$
- 4) train a bias LM with the selected passages
- 5) estimate the interpolation weight, $\alpha(d)$
- 6) compute log LM scores according to (Eqn. 1)

The CLIR tool we used in Step 1 is the one presented in (Xu, etc., 2001), where details of the CLIR score computation can be found. In Step 3, “ T ” is a threshold that controls the selection. In Step 4, when estimating the bias LM, we applied higher weighting on n-grams counted from more closely related passages. The weighting factor for n-gram counts from a selected passage is computed according to

$$\text{psg_wgt} = \frac{\text{CLIR_scr}(p) - \mu(d)}{\arg \max_{\{p \in \text{selected_psgs}\}} \text{CLIR_scr}(p) - \mu(d)}$$

where “ $\text{CLIR_scr}(p)$ ” represents the CLIR score of a selected passage.

In Step 5 we estimated the interpolation weight according to

$$\alpha(d) = \frac{\mu_s(d) - \mu(d)}{\arg \max_{\{d' \in \text{all_test_docs}\}} \mu_s(d') - \mu(d')}$$

where $\mu_s(d)$ represent the means of the CLIR scores of all the selected passages. The weight is normalized to be between 0 and 1.

We used our own CLIR tool to extract related passages. We found that it was a good choice to set $N = 4,000$ in Step 1 and set the selection threshold “ T ” to 1.28 in Step 3. The related passages for training the bias LM were extracted from the 14B English patent document corpus. In our case here passages are equivalent to patent documents in the English patent corpus.

System	Test
<i>Add top-100 best features + 14B-pt LM</i>	39.14
<i>+ LM adaptation</i>	40.04
<i>LM adaptation (patent description)</i>	39.97
<i>LM adaptation (patent abstract)</i>	40.23

Table 5. Improvements (in terms of the BLEU score) from the LM adaptation

To have the right baseline system, we re-ran the “*add top-100 best features*” system, shown in Table 4, but switching from the “45M-pt” LM to the “14B-pt” LM. This system is denoted as “*Add top-100 best features + 14B-pt LM*” in Table 5. Comparing these two systems, we see that the “14B-pt” LM improved the MT performance by 1.4 BLEU points (from 37.71 to 39.14). We then conducted the LM adaptation on top of the “*Add top-100 best features + 14B-pt LM*” system. The performance is shown in the “*+ LM adaptation*” row in Table 5. The LM adaptation improved the BLEU score by 0.9 points (from 39.14 to 40.04).

As described earlier, the 2K sentences of the development set were extracted from the descriptions of 103 patents, so the patent documents in the Test set are only portions of the corresponding original patent documents. In the above LM adaptation we used the portions of the full patent descriptions as queries for the CLIR. Since the NTCIR-9 organizers also provided the full contents of the 103 patent documents, we explored uses of the abstracts and the full descriptions of the patent documents as queries for the CLIR in the LM adaptation. The scores of these two experiments are listed in the last two rows in Table 5. As shown, the uses of the abstracts and full descriptions in the LM adaptation produced similar results.

5 Conclusion

We have described the work we carried out for building an SMT system for the Chinese-English patent MT sub-task of the NTCIR-9 MT evaluation. First, we made changes to our SMT training procedure in order to better handle the special characteristics of patent data, and obtained incremental improvements from the various changes. Then, the re-training of the casing LM with patent text and the use of more features to the MT system improved the BLEU scores significantly. Finally, the LM adaptation improved the MT performance further – by about 1 BLEU point. Our work shows that most of the strategies for building an SMT system, such as the use of a large number of features and the LM adaptation, were easily applicable to the patent genre and produced gains. But certain techniques had to be customized in order to better handle the special characteristics of the patent genre.

References

- Stanley F. Chen and Joshua Goodman, 1998. “An Empirical Study of Smoothing Techniques for Language Modeling”, *Tech. Rep. TR-10-98, Harvard University*.
- David Chiang, Kevin Knight, and Wei Wang. 2009. “11,001 new features for statistical machine translation”, in *Proceedings of the Human Language Technologies, NAACL 2009*.
- Makoto Iwayama, Fujii Atsushi, Kando Noriko, and Akihiko Takano. 2003. “Overview of patent retrieval task at NTCIR-3”, in *Proceedings of the third NTCIR workshop*.
- Yaohong Jin. 2010. “A hybrid-strategy method combining semantic analysis with rule-based MT for patent machine translation”, in *Proceedings of 2010 International Conference on NLP-KE, 2010*.
- Jae-Ho Kim and Key-Sun Choi. 2007. “Patent document categorization based on semantic structural information”, *Information Processing and Management (2007)*, doi:10.1016/j.ipm.2007.02.002
- Franz J. Och and Hermann Ney. 2003. “A systematic comparison of various statistical alignment models”, *Computational Linguistics*, 29(1):19_51.
- Len Offersgaard and Claus Povlsen. 2007. “Patent documentation – comparison of two MT strategies”, in *Proceedings of the Machine Translation Summit XI workshop on Patent Translation*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. “BLEU: a Method for Automatic Evaluation of Machine Translation”, in *Proceedings of the 40th Annual conference of the association of computational linguistics*.
- Atti-Veikko Rosti, Bing Zhang, Spyros Matsoukas, and Richard Schwartz. 2010. “BBN System Description for WMT 10 System Combination Task”, *Proceedings of the 5th workshop on Statistical Machine Translation, 2010*.
- Sayori Shimohata. 2005. “Finding Translation Candidates from Patent Corpus,” MT Summit X Workshop on Patent Translation, 2005.
- Libin Shen, Jinxi Xu, and Ralph Weischedel. 2008. “A new string-to-dependency machine translation algorithm with a target dependency language model”, *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics*.
- Matthew Snover, Bonnie J. Dorr, and Richard Schwartz. 2008. “Language and translation model adaptation using comparable corpora”, in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- Matthew Snover, Nitin Madnani, Bonnie Dorr, and Richard Schwartz. 2010. “TER-Plus: Paraphrase, Semantic, and Alignment Enhancements to Translation Edit Rate”, *Machine Translation Journal, Special Issue on: Automated Metrics for MT Evaluation*.
- Terumasa Ehara. 2007. “Rule Based Machine Translation Combined with Statistical Post Editor for Japanese to English Patent Translation”, in *Proceeding of MT Summit XI*, pp.13-18.
- Dan Wang. 2009. “Chinese to English automatic patent machine translation at SIPO”, *World Patent Information, Volume 31, Issue 2*, pp. 137-139.
- Jinxi Xu, Ralph Weischedel and Chanh Nguyen. 2001. “Evaluating a Probabilistic Model for Cross-Lingual Information Retrieval”, in *SIGIR 2001: 105-110*