

Learning of Linear Ordering Problems and its Application to J-E Patent Translation in NTCIR-9 PatentMT

Shuhei Kondo
 Nara Institute of Science and Technology
 8916-5 Takayama, Ikoma, Nara 630-0192, Japan
 shuhei-k@is.naist.jp

Katsuhiro Sudoh
 NTT Communication Science Laboratories
 2-4 Hikaridai, Seika-cho, Soraku-gun, Kyoto 619-0237 Japan
 sudoh.katsuhiro@lab.ntt.co.jp

Mamoru Komachi
 Nara Institute of Science and Technology
 8916-5 Takayama, Ikoma, Nara 630-0192, Japan
 komachi@is.naist.jp

Kevin Duh
 NTT Communication Science Laboratories
 2-4 Hikaridai, Seika-cho, Soraku-gun, Kyoto 619-0237 Japan
 kevin.duh@lab.ntt.co.jp

Yuji Matsumoto
 Nara Institute of Science and Technology
 8916-5 Takayama, Ikoma, Nara 630-0192, Japan
 matsu@is.naist.jp

Hajime Tsukada
 NTT Communication Science Laboratories
 2-4 Hikaridai, Seika-cho, Soraku-gun, Kyoto 619-0237 Japan
 tsukada.hajime@lab.ntt.co.jp

ABSTRACT

This paper describes the patent translation system submitted for the NTCIR-9 PatentMT task. We applied the Linear Ordering Problem (LOP) based reordering model [16] to Japanese-to-English translation to deal with the substantial difference in the word order between the two languages.

Categories and Subject Descriptors

I.2.7 [Artificial Intelligence]: Natural Language Processing—*machine translation*

General Terms

Experimentation

Keywords

Statistical Machine Translation, Reordering, Preprocessing.

Team Name

[NAIST]

Subtasks/Languages

[Japanese to English]

External Resources Used

[Moses][GIZA++][SRILM][MeCab]

1. INTRODUCTION

While phrase-based statistical machine translation can model short-distance reordering within the phrase fairly well, long-distance reordering remains one of the major problems in SMT, especially for translations between language pairs with substantially different word order.

Japanese and English are an example of such a pair, the former being an SOV language and the latter an SVO language. Translations between such language pairs are difficult

because the decoder often misses long-distance reordering required for an appropriate translation, partly due to the distortion limit that is commonly adopted by phrase-based decoders to reduce the search space.

Reordering by preprocessing is one of the common approaches to deal with long-distance reordering required for translation between languages with different word order. Works in this direction can be classified into syntax-based [3, 9], chunk-based [1, 18] and word-based [9, 16] approaches. The syntax-based and the chunk-based approaches depend on syntactic parsing and base-phrase chunking, and therefore they usually require rich language resources such as treebanks. In contrast, word-based approaches only use part-of-speech (POS) tagger or simple rules, and thus it is applicable to languages and domains with limited resources.

We applied one of the word-based approaches, the LOP-based word-level reordering model [16], for Japanese-to-English translation in the system submitted to the NTCIR-9 patent machine translation task [7].

2. LINEAR ORDERING PROBLEM BASED REORDERING

Tromble and Eisner (2009) [16] proposed a word-level reordering model based on Linear Ordering Problem (LOP) and an efficient decoding algorithm called “local search”.

2.1 The LOP-based Model

Given an input sentence $\mathbf{w} = w_1 w_2 \dots w_n$, a permutation $\pi = \pi_1 \pi_2 \dots \pi_n$ is any reordering of the tokens in \mathbf{w} . The LOP-based model is defined so that it assigns a high score to the permutation π that corresponds well to the word order of the target language. To construct such model, a pairwise preference matrix $B_{\mathbf{w}}$ is constructed for each sentence \mathbf{w} , whose entries are

$$B_{\mathbf{w}}[l, r] \stackrel{\text{def}}{=} \theta \cdot \phi(\mathbf{w}, l, r) \quad (1)$$

where θ is a vector of weights, ϕ is a vector of feature functions considering the input sentence \mathbf{w} and its functions such as POS tags, l and r are the positions of left and right

Table 1: Feature templates for $B_w[l, r]$ (w_l is the l th word and t_l is its POS tag. b matches any index with $l < b < r$.) Each combination is conjoined with the distance between the words, $r - l$. Distances are binned into 1, 2, 3, 4, 5, $5 < d \leq 10$ and $10 < d$.

| t_{l-1} | w_l | t_l | t_{l+1} | t_b | t_{r-1} | w_r | t_r | t_{r+1} |
|-----------|-------|-------|-----------|-------|-----------|-------|-------|-----------|
| | × | × | | | | × | × | |
| | × | × | | | | × | | |
| | × | | | | × | × | | |
| | × | × | | | | × | | |
| | | × | | | × | × | | |
| | | × | | | | × | | |
| | | × | | | × | | | |
| | | × | | | | × | | |
| | | × | | | | | × | |
| | | | × | | | | × | |
| | | | × | | | | | × |
| | | | | × | | | | |
| | | | | | × | | | |
| | | | | | | × | | |
| | | | | | | | × | |
| | | | | | | | | × |
| × | | | | | | | | |
| × | | | | | | | | |
| × | | | | | | | | |
| × | | | | | | | | |

word in the input sentence. $B\mathbf{w}[l, r] > B\mathbf{w}[r, l]$ means that the model prefers a permutation in which the order of the l -th word and r -th word is the same as the input, while $B\mathbf{w}[l, r] < B\mathbf{w}[r, l]$ means that the order of l -th word and r -th word should be reversed. The score for a permutation π is defined as

$$\text{score}(\pi) \stackrel{\text{def}}{=} \sum_{i, j: 1 \leq i < j \leq n} B\mathbf{w}[\pi_i, \pi_j], \quad (2)$$

where $B\mathbf{w}[r, l]$ can be considered 0 for any $l < r$, because subtracting $B\mathbf{w}[r, l]$ from both $B\mathbf{w}[l, r]$ and $B\mathbf{w}[r, l]$ will reduce the scores of all permutations by the same amount.

To calculate the score for (2), the weight vector θ in (1) is learned discriminatively with the reference reordering obtained from automatic word alignments. Each source token is assigned an integer key that corresponds to the position of the leftmost target token which is aligned to it, or 0 if it is not aligned to any target token. The order of l -th token and the r -th token is swapped in the reference reordering if the key for the former is larger than the latter. Table 1 shows the feature templates for ϕ used by [16], which is adapted from [11].

2.2 Local Search

An algorithm called “local search” is used to reorder the source sentence according to the LOP-based model. Given an input sentence \mathbf{w} , for each span $i \dots k$ which yields the substring $w_{i+1} w_{i+2} \dots w_k$, the score $\max(0, \Delta_{i,j,k})$ for its subspan $i \dots j, j \dots k$ is calculated according to the model for all pairs of subspans. $\Delta_{i,j,k}$ is the score of swapping the subspans with the base score $\Delta_{i,j,k} = 0$ for $i = j$ and $j = k$, and a positive score means that they should be swapped. Starting with

Table 2: BLEU scores on the development set with varying distortion limits.

| Word Order | Distortion Limit | BLEU |
|------------|------------------|-------|
| LOP-based | 6 | 29.52 |
| LOP-based | 10 | 29.71 |
| LOP-based | 20 | 29.56 |
| Original | 10 | 28.56 |
| Original | 20 | 29.44 |
| Original | 30 | 29.30 |

Table 3: Average execution time (in seconds) of the LOP-based reordering of 10,000 sentences in 10 trials. (a) is for the first 10,000 sentences in the training set, whose average length is 19.60 words, and (b) is for 10,000 sentences from the 1,000,001st to 1,010,000th sentence, whose average length is 36.29 words.

| | (a) | (b) |
|-----------|-------|--------|
| time(sec) | 56.53 | 337.23 |

Table 4: Official automatic evaluation results

| System | BLEU | NIST | RIBES |
|-----------|-------|--------|----------|
| LOP-based | 27.82 | 7.4348 | 0.730743 |
| BASELINE1 | 28.95 | 7.7696 | 0.70644 |
| BASELINE2 | 28.61 | 7.7562 | 0.675831 |

Table 5: Official human evaluation results

| System | Adequacy | Acceptability |
|-----------|----------|---------------|
| LOP-based | 2.610 | 0.472 |
| BASELINE1 | 2.617 | 0.474 |
| BASELINE2 | 2.427 | 0.447 |

Table 6: BLEU scores on the NTCIR-9 PatentMT test set

| Word Order | BLEU | NIST | RIBES |
|------------|-------|--------|----------|
| LOP-based | 27.82 | 7.4348 | 0.730737 |
| Original | 27.37 | 7.5648 | 0.675202 |

the spans that consist of one token and recursively expanding the spans, local search finds the best permutation of the input sentence under the ITG [17] constraints according to the LOP-based model in $O(n^3)$ time, where n is the length of the input sentence.

3 EXPERIMENTS

3.1 Data

We used the training set of NTCIR-7 Patent Translation Task (1.8 million sentence pairs) [5, 6] to train the translation model, and its English part for the language model; data from 2001 to 2005 were not used due to time limitation. Sentences with more than 80 tokens were also excluded. We corrected some obvious errors in its Japanese part, such as prolonged sound marks replaced by minus signs, using regular expressions.

To learn the LOP-based reordering model, we picked 300,000 sentence pairs with high sentence alignment score from the training set. The weight vector θ was learned using Averaged Perceptron [2] with the 300,000 pairs. We limited the amount of corpus here because automatic word alignments obtained from sentence pairs with low alignment score are

likely to be less reliable and they can be harmful to the reordering model.

Then training, development and test sentences were reordered by local search with the learned weights, and the translation model was retrained on the reordered training set. Table 2 shows the BLEU [14] scores of the LOP-based system on the development set with varying distortion limits, compared with the system with original word order and same amount of corpus. The result shows that the LOP-based system requires smaller distortion limit. As the LOP-based system performed best with $d = 10$ on the development set, we used this value for the test.

Table 3 shows the average execution time required for reordering of source sentences. We can see reordering by preprocessing in our system can be performed quite efficiently, even with the naive implementation written in Ruby for reordering.

3.2 Tools

Japanese sentences were tokenized and POS-tagged using MeCab version 0.98¹. We used mecab-ipadic-2.7.0-20070801 as the dictionary, with a slight modification on costs for the words whose POS tag is exclamation, to penalize the words that are unlikely to appear in patent texts.

We used GIZA++ (version 1.0.5) [13] for word alignment, SRILM (version 1.5.11) [15] for construction of 5-gram language model and Moses (revision 3947) [10] for training of translation models and decoding. Parameters for each model were tuned using minimum error-rate training [12] with BLEU as the objective. Recaser was trained in similar manner to organizer's baseline with the same amount of data as the training set of our system.

3.3 Results

Table 4 shows the official scores of our system based on automatic evaluation metrics, compared with baseline systems provided by the organizer. Our system achieved lower scores in BLEU and NIST [4] metric, and higher score in RIBES [8] metric than both baseline systems.

Table 5 shows the official scores of our system based on human evaluations. Our system's performance was comparable to BASELINE1 (hierarchical phrase based) and slightly better than BASELINE2 (phrase-based). This result seems encouraging, because our system was trained on less resources and its procedure after preprocessing was quite similar to that of BASELINE2.

4. ANALYSIS

We compared the LOP-based system with a system trained with the same amount of corpus that is in the original word order, because the LOP-based system was trained on smaller corpus compared to the BASELINE systems provided by the organizer as mentioned above in section 3.

Table 6 shows the BLEU, NIST and RIBES scores of both systems. While the reordered system's performance was worse in the NIST score, it performed better in BLEU and RIBES scores. Both in this result and the result in the official test, the LOP-based model seems to perform well in RIBES score and poorly in NIST score.

We looked at example outputs of both systems to analyze the reason for this tendency. Table 7 shows several examples

¹<http://mecab.sourceforge.net/>

of original source sentences, reordered inputs, outputs of the system trained with original word order, and outputs of the LOP-based system with references, in which the LOP-based reordering model worked well. In these examples, the reordered inputs can be translated almost monotonically.

However, we can also see that some case-markers are moved out of their context in the reordered inputs. Moreover, in the last example, tokens forming the Japanese compound word for "engine control system" are separated as a result of reordering. Though the effects on overall reordering were small in these examples, such displacement can harm the quality of translation.

The gap between BLEU and NIST scores can possibly be attributed to this displacement of case-markers and compound words, because NIST metric gives more weight to rare n-grams and therefore is more sensitive than BLEU metric to such displacement, especially in the cases with compound words. One way to deal with this problem might be to change the unit of reordering, for example from words to chunks. Extending the unit of reordering to ones containing more than one words might contribute to reduction of the time required for reordering, as local search is $O(n^3)$ algorithm and depends on the number of units to be reordered.

High RIBES score may indicate that overall reordering performed by our system is reasonably good, because RIBES metric is designed so as to evaluate the global word reordering well, which is required for translation with distant language pairs. Parameter tuning with RIBES score instead of BLEU score as the objective might further improve the global reordering.

5. CONCLUSION

In this paper, we applied the LOP-based word-level reordering for Japanese-to-English translation. In automatic evaluation, our system performed worse in two metrics and better in one metric compared to both baseline systems. When compared with the same amount of training data, our system performed better in two metrics than the system with original word order. In human evaluation, our system performed better than phrase-based baseline and similarly to hierarchical phrase-based one, despite the fact that our system was trained on less amount of data.

Besides the problem of global reordering, displacement of case-markers and tokens that form compound words was observed in our system. We think one way of the future work will be to find the model that can prevent such displacement.

6. ACKNOWLEDGMENTS

The authors would like to thank Katsuhiko Hayashi for valuable comments and discussions, and the anonymous reviewer for insightful comments and suggestions.

7. REFERENCES

- [1] A. Bisazza and M. Federico. Chunk-based verb reordering in VSO sentences for Arabic-English statistical machine translation. In *Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and MetricsMATR*, pages 235–243, 2010.
- [2] M. Collins. Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms. In *Proceedings of the 2002*

Table 7: Sample of reordered inputs and outputs of our system

| | |
|--------------------------|--|
| Original source sentence | そして、スピナバルブ型巨大磁気抵抗効果素子対2, 3, 4, 5 それぞれの他方の端部間に測定用電圧Vccが印加されている。 |
| Reordered input | そして、のにてが印加されている測定用電圧Vcc間他方端部 のそれぞれスピナバルブ型巨大磁気抵抗効果素子対2, 3, 4, 5。 |
| Original output | Then, the spin valve type giant magnetoresistance effect element pair 2, 3, 4, 5, respectively, between the other end of the measuring voltage Vcc is applied. |
| LOP-based output | Then, the measuring voltage Vcc is applied to between the other end of the spin-valve type giant magnetoresistive element pair 2, 3, 4, 5. |
| Reference | A voltage Vcc for measurement is applied between the other ends of the spin-valve giant magneto-resistive element pairs 2, 3, 4 and 5. |
| Original source sentence | 横方向スキャン装置70の動作は、レーダ装置10の制御演算部11にて制御される。 |
| Reordered input | 動作の横方向スキャン装置70は制御されるにての制御演算部11、レーダ装置10。 |
| Original output | The operation of the radar apparatus 10, the control arithmetic unit 11 in the lateral direction is controlled by a scanning device 70. |
| LOP-based output | The operation of the lateral direction scan unit 70 is controlled by the control arithmetic unit 11 of the radar apparatus 10. |
| Reference | The operations of the horizontal-direction scanning device 70 are controlled by the control and signal processing unit 11 of the radar apparatus 10. |
| Original source sentence | この第3過程における位相変調は、集光レンズ等からなる 第2の光学素子12を用いて実現できる。 |
| Reordered input | この位相変調における第3過程はできる実現を用いての 第2光学素子12からなる、集光レンズ等。 |
| Original output | In the third process, a condenser lens or the like using the phase modulation of the second optical element 12 can be realized. |
| LOP-based output | This phase modulation in the third process is implemented using a second optical element 12 comprising the condensing lens or the like. |
| Reference | The phase modulation in this third process can be realized using a second optical element 12 composed of a focusing lens or the like. |
| Original source sentence | マップデータCは、センサ故障などエンジン制御システムに 異常が発生した場合に使用する特性変換係数用のマップデータである。 |
| Reordered input | マップデータCはでのマップデータ用ある特性変換する係数使用 に場合異常が発生したに、エンジンなどセンサ故障制御システム。 |
| Original output | Map data C, the sensor failure or abnormality has occurred in the engine control system used in a case where a characteristics conversion coefficient map data. |
| LOP-based output | The map data C is a map data for the characteristic conversion coefficient used in a case where abnormality is generated in the engine sensor failure control system. |
| Reference | The map data C is the map data for the characteristic conversion coefficient used when a defect occurs in the engine control system, such as a failure of a sensor. |

- Conference on Empirical Methods in Natural Language Processing, pages 1–8, 2002.
- [3] M. Collins, P. Koehn, and I. Kučerová. Clause restructuring for statistical machine translation. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 531–540, 2005.
- [4] G. Doddington. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of the second international conference on Human Language Technology Research*, pages 138–145, 2002.
- [5] A. Fujii, M. Utiyama, M. Yamamoto, and T. Utsuro. Overview of the Patent Translation Task at the NTCIR-7 Workshop. In *Proceedings of NTCIR-7 Workshop Meeting*, pages 389–400, 2008.
- [6] A. Fujii, M. Utiyama, M. Yamamoto, T. Utsuro, T. Ehara, H. Echizen-ya, and S. Shimohata. Overview of the Patent Translation Task at the NTCIR-8 Workshop. In *Proceedings of NTCIR-8 Workshop Meeting*, pages 371–376, 2010.
- [7] I. Goto, B. Lu, K. P. Chow, E. Sumita, and B. K. Tsou. Overview of the Patent Machine Translation Task at the NTCIR-9 Workshop. In *Proceedings of NTCIR-9 Workshop Meeting*, 2011.
- [8] H. Isozaki, T. Hirao, K. Duh, K. Sudoh, and H. Tsukada. Automatic Evaluation of Translation Quality for Distant Language Pairs. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 944–952, 2010.
- [9] J. Katz-Brown and M. Collins. Syntactic reordering in preprocessing for Japanese → English translation: MIT system description for NTCIR-7 patent translation task. In *Proceedings of NTCIR-7 Workshop Meeting*, pages 409–414, 2008.
- [10] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster*

- and Demonstration Sessions*, pages 177–180, 2007.
- [11] R. McDonald, K. Crammer, and F. Pereira. Spanning Tree Methods for Discriminative Training of Dependency Parsers. Technical Report MS-CIS-05-11, UPenn CIS, 2005.
 - [12] F. J. Och. Minimum error rate training in statistical machine translation. In *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics*, volume 1, pages 160–167, 2003.
 - [13] F. J. Och and H. Ney. A systematic comparison of various statistical alignment models. *Computational linguistics*, 29(1):19–51, 2003.
 - [14] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
 - [15] A. Stolcke. SRILM—an extensible language modeling toolkit. In *Proceedings of the International Conference on Spoken Language Processing*, volume 2, pages 901–904, 2002.
 - [16] R. Tromble and J. Eisner. Learning Linear Ordering Problems for Better Translation. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 1007–1016, 2009.
 - [17] D. Wu. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*, 23(3):377–403, 1997.
 - [18] Y. Zhang, R. Zens, and H. Ney. Improved chunk-level reordering for statistical machine translation. In *Proceedings of the International Workshop on Spoken Language Translation*, pages 21–28, 2007.