# Improved Reordering Rules for Hierarchical Phrase-based Translation

Shu Cai, Yajuan Lü, Qun Liu

*Key Lab. of Intelligent Information Processing*
*Institute of Computing Technology, Chinese Academy of Sciences*
*P.O. Box 2704, Beijing 100190, China*
{*caishu,lvyajuan,liuqun*}*@ict.ac.cn*

*Abstract*—**Hierarchical phrase-based translation model has been proven to be a simple and powerful machine translation model. However, due to the computational complexity constraints, the extraction and use of hierarchical rules are usually restricted under certain limits, and these limits could have a negative impact on the performance of the translation model, especially for reordering. This paper presents a solution to improve the reordering of hierarchical phrase-based translation model. We propose a two-step method to extract** *improved reordering rules* **with less limits. These reordering rules help both local and non-local reordering, and could be incorporated to a hierarchical phrase-based translation system easily. Experiments show that our approach achieves statistically significant improvements over the baseline system in Chinese-English translation.**

*Keywords*-**statistical machine translation; hierarchical phrase-based translation; hierarchical rules; reordering;**

## I. INTRODUCTION

Hierarchical phrase-based translation model [1] combines the ideas of syntax-based translation and phrase-based translation. This model extracts rules of synchronous context-free grammar (*hierarchical rules*) from a word-aligned corpus and match these rules hierarchically during decoding. The implementation of this model, the Hiero system[1], has been proven to be powerful in the machine translation community.

Reordering problem is a central yet difficult problem in machine translation. Generally, there are two kinds of approaches for reordering: the first is reordering the input sentence before translation, such as [2]; the second is adding a reordering model as an additional feature, such as [3], [4]. Different from all these approaches, hierarchical phrase-based translation implicitly solves reordering problem during hierarchical rule matching. This solution is simple and effective since reordering information is encoded in the rules and we do not have to calculate reordering probabilities for each span.

However, due to the computational complexity, the extraction and use of hierarchical rules are usually imposed with many constraints, including constraints on the length of initial phrases, number of nonterminals plus terminals, adjacent nonterminals .etc. These constraints may successfully ensure the efficiency of extraction and decoding, but often fail to capture useful patterns, especially for translation between language pairs with much non-local reordering, such as Chinese and English.

Several recent works attempt to improve the reordering of hierarchical translation model. [5] combines shallow tree-to-string rules with hierarchical rules and uses them to help reordering of the hierarchical phrase-based translation model. This method takes advantage of the linguistic analysis results, and introduce a great number of additional rules to the existing translation system.

In this paper, we proposes a light-weight method to improve the reordering of hierarchical phrase-based translation systems. Different from the extraction in [1], our extraction method only considers hierarchial rules which help reordering, and uses a two-step process to extract these rules and calculate their probabilities. These rules has less constraints than rules in [1]. Different from [5], the amount of our rules is small and we do not use parser of the source language. In addition, it is simple and effective to incorporate these new rules into a hierarchical phrase-based translation system by slight modification. Experiments show that after adding these rules, the translation quality of the system improves significantly in Chinese-to-English translation.

The rest of the paper is organized as follows: rules used in hierarchical phrase-based translation model is briefly introduced in Section II. The extraction of the improved reordering rules are described in Section III. In Section IV, we describe the decoding procedure which incorporated reordering rules. Experiments' details are shown in Section V, and the conclusions are shown in the Section VI.

## II. HIERARCHICAL RULES

In [1], the hierarchical rules have the following form:

- $X \rightarrow < \gamma, \alpha, \sim >$

where X is a nonterminal, $\gamma$ and $\alpha$ are both strings of terminals and nonterminals, $\sim$ is a one-to-one correspondence between nonterminal occurrences in $\gamma$ and $\alpha$ and is often omitted when writing the rules.

Some Chinese-English hierarchical rule examples could be written as follows:

- $X \rightarrow < 今天 \ X1, \ X1 \ today >$
- $X \rightarrow < X1 \ 的 \ X2 , \ X2 \ of \ X1 >$

The nonterminals of the rules are generated by replacing subphrases in the initial phrases with variables. Initial

phrases are extracted from the corpus according to word-alignment consistency in [6].

To reduce the rule set size and spurious ambiguity, the following constraints are used to filter the rules:

(1) Unaligned words are not allowed at the edges of initial phrases.

(2) Initial phrases are limited to a certain length on either side. (10 words in [1])

(3) The number of nonterminals plus terminals are restricted to a certain number (5 in [1]).

(4) Rules can have at most two nonterminals.

(5) Adjacent nonterminals are not allowed on the source side.

(6) The rule must have at least one pair of aligned words.

Each rule is associated with 4 probabilities: frequency probabilities $P(\gamma|\alpha)$, $P(\alpha|\gamma)$, lexical frequency probabilities $P_w(\gamma|\alpha)$, $P_w(\alpha|\gamma)$.

We refer to hierarchical rules in [1] as the *ordinary (hierarchical) rules* afterwards.

In decoding, glue rules

- S → < S X, S X >
- S → < X , X >

are used when no rules could match or the span exceeds a certain length (*search depth*). Glue rules simply monotonically connect translations of two adjacent blocks together. The search depth is often set the same as the initial phrase length limit.

## III. IMPROVED REORDERING RULE EXTRACTION

Since the extraction of ordinary hierarchical rules may replace any subphrase with variables, the resulting rule set' s size could be huge. Adding constraints may ensure the efficiency of extraction and decoding, but will also filter useful patterns. Although in theory, we may extract ordinary hierarchical rules without any constraints, the resulting rule set will be too large to use. For example, if we change the initial length limit from 10 to 15 in ordinary rule extraction of FBIS corpus, about 4 million rules will be added.

In decoding, glue rule is often used as the only rule to connect two spans when there are no matching hierarchical rules. However, this solution neglects the reordering possibilities and may degrade the performance. The *improved reordering rules* are extracted only from initial phrases containing reordering examples in the corpus. These rules try to capture useful reordering patterns ignored by using various constraints in the ordinary rule extraction. In the following section, we present three kinds of improved reordering rules with different constraints relaxed, and show the extraction process.

### A. Improved Reordering Rule

First, we extract the improved reordering rules without length limit on initial phrase length or number of nonterminals plus terminals. Patterns which cover a long range

might be ignored by using fixed limit in hierarchical rule extraction. In the following example, if we set the initial phrase limit to 10 words (as in [1]), we will fail to extract (2) from (1) although (2) is a reasonable rule to describe the structure of the sentence pair (1).

(1) 澳洲 是 与 北 韩 有 邦 交 且 与 南 韩 也 有 良 好 关 系 的 少 数 国 家 之一

(Australia is with North Korea have diplomatic relationship and with South Korea also have good relations that few countries one of )

"Australia is one of the few countries that have diplomatic relations with North Korea and have good relations with South Korea"

(2) 是 X1 之一

(is X1 one of)

is one of X1

(2) is a reordering pattern which perform a long-distance reordering. Without it, simply using glue rule in decoding will translate the Chinese part of (1) incorrectly.

Statistics show that such long reordering examples are quite common in some language pairs. In the FBIS corpus, which consists of 239,371 Chinese-English sentence pairs, there is 141,299 reordering examples whose source side is longer than 10 words. In the ordinary rule extraction, rules extracted from these reordering examples are pruned.

Second, improved reordering rules allow rules with consecutive variables on the source side. The ordinary hierarchical rules forbid consecutive variables on the source side since they are main causes of ambiguity. But the consecutive variables play an important role in reordering patterns. To illustrate,

- X → < 在 X1 X2 , X2 in X1 >

is a frequent reordering pattern in Chinese-English when X1 is some country and X2 is some activity.

Rules with consecutive variables are also common in reordering patterns. The above example occurs 268 times in the FBIS corpus, and the rules with consecutive variables has a total count of 27070.

Finally, improved reordering rules allow more than two nonterminals. To provide more context for the reordering, we extract some reordering patterns by analyzing the hierarchical structure of the language. Since initial phrases containing reordering examples may contain other such phrases, we replace part of the reordering rule extracted from longer phrases with the reordering rule extracted from shorter phrases to get new reordering rules. It might result in rules with more than two variables.

This process is weakly equivalent to the decoding process of hierarchical phrase-based translation. Although in theory, hierarchically matching the ordinary hierarchical rules is enough for generating translations, limited context in matching reordering rules in a long span is prone to errors. In the following example,

(1) X → < X1 所 X2 , X2 by X1 >
(2) X → < X1 的 X2 , X2 of X1 >
(3) X → < X1 所 X2 的 X3 , X3 X2 in the X1 >

where (1) (2) are the most frequent rules used in the limited context. But when we combine the context together in (3), the most frequent rule is not the combination of the two most frequent rules above. We could see that rules which describe the hierarchical structure provide a better guidance for reordering during the translation.

In all, the improved reordering rules capture the reordering patterns by analyzing all reordering examples. By relaxing various limits, we could capture complete reordering patterns. Not only could these patterns be used in both long span and short span during decoding in hierarchical phrase-based translation, but also could they be used in other phrase-based translation to guide the reordering.

To reduce spurious ambiguity and gain efficiency, we use Part-Of-Speech(POS) tags of the source language to help decide which part of the reordering rules should be replaced by variables. We require that boundary words (the head word and the tail word) of variables are both content words[1]. POS tags and boundary words are shown to be useful features of reordering in [8], [9], [10]. The POS tag information is associated with the improved reordering rules to be used in decoding.

Improved reordering rules have the same form as the ordinary hierarchical rules, along with the POS tags of variables' boundary words. An example of improved reordering rules is as follows:

- X → < X1 的 X2 , X2 of X1 , CD NN JJ NN >

where "X1 的 X2" is the source string, "X2 of X1" is the target string, "CD NN" are POS tags of the boundary words of X1 (on the source side) , and "JJ NN" are those of X2.

In the improved reordering rule extraction, we keep the constraints (1), (6) of the ordinary hierarchical rules (Section II), and relax the other constraints. We further filter some reordering rules according to its pattern because these patterns have too many consecutive non-terminals and too little lexical evidence, which may easily cause ambiguities. These rules include rules which have more than 2 consecutive variables or more than 3 variables, and rules whose lexical parts are all punctuation.

### B. Two-step Extraction

Since we only extract rules from reordering examples, the time complexity is lower compared to the ordinary rule extraction. However, we do not have the whole rule set to estimate the frequency probabilities. To solve this problem, we use a two-step extraction process. First, we extract the reordering rules and calculate their lexical probabilities. Second, we get the counts of each rule's source and target

---

[1]In the tag set of Chinese Penn Treebank [7], we define the tag set { *NN, JJ, VV, NR, CD, DT, PN*} to be the content word tags.

side by matching improved reordering rules in the corpus, then calculate the rules' frequency probabilities.

Given a bilingual sentence, the first step of our rule extraction algorithm is to extract improved reordering rules from all initial phrases containing reordering examples. The source side of these initial phrase pairs could generally be viewed as two adjacent blocks whose translation is still adjacent, but reordered on the target side. In our extraction, the source part is first divided into such adjacent blocks. For each block, we search for the longest part whose boundary words are both content words in each block and replace this part with a variable. A special case is that there might be more than one way to divide the phrase. Since the cause of this case is that a "middle word" exists between the two reordered blocks and this word could be divided into any of the blocks, we divide the word to the right block. Actually either division results in the same rule.

Each adjacent block of an initial phrase pair may contain other reordered phrases. Thus, we could build a tree structure of reordering phrases to show this hierarchical structure. By using a depth-first search, we analyze the reordering patterns from shortest phrases to longer phrases which contain shorter phrases. If a phrase contains shorter phrases, we replace parts in reordering rules extracted from it with shorter reordering rules to get new rules.

Figure 1 (a) shows the word alignments of a phrase pair. The original phrase is quite long, and we only show the main part of it for simplicity. The extraction is briefly shown in Figure 1 (b). The tree structures of the phrase's source part and target are shown in the upper half and lower half of the Figure respectively, linked by dotted arrows. The phrase's source part is labeled with POS tags. The dashed lines in Figure 1 (b) show the variable-generating process for each reordered block, and the lines between variable pieces indicate the correspondences. The improved reordering rules extracted from this phrase are shown in Table I in order. We search the tree from bottom to top and extract these rules. The first rule is extracted from "邦交 的 政策" and its translations. (Note that the source part could be divided in two ways and results in the same rule. ) The second rule is extracted from the longer phrase "和谈 后 确定 邦交 的 政策" and its translation. The third rule is obtained by replacing the variable in the second rule (extracted from longer part of the phrase) with the first rule (extracted from shorter part of the phrase), and the part other than the reordering part "确定" and its translation "decide" are kept.

It is easy to calculate the two lexical probabilities $P_w(\gamma|\alpha)$, $P_w(\alpha|\gamma)$ of the improved reordering rules according to the alignments. But translation probabilities $P(\gamma|\alpha)$, $P(\alpha|\gamma)$ could not be estimated directly after the first step. Thus in the second step of the extraction, we go through the corpus again, trying to match improved reordering rules in the sentences. After we get the total occurrence number of source side and target side for each improved reordering rule,
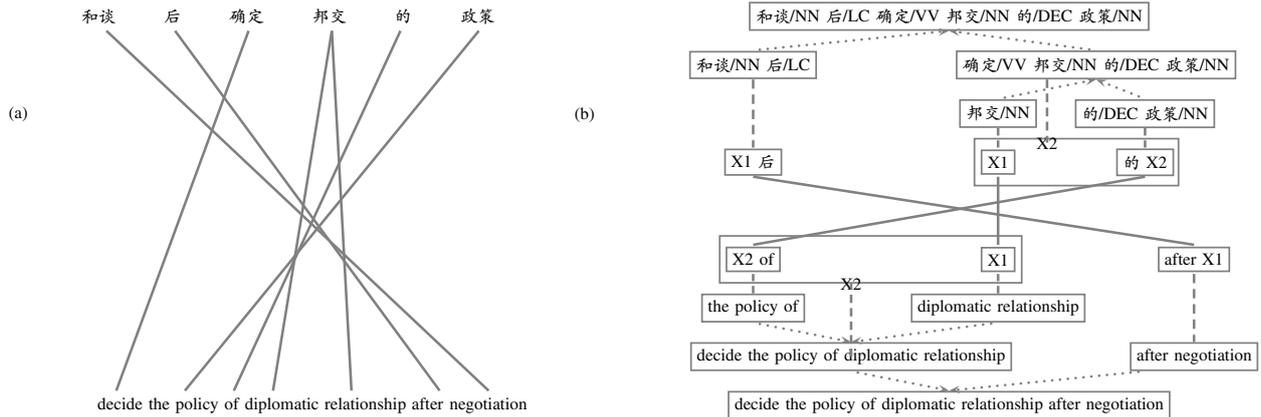
Figure 1. Reordering rule extraction of an initial phrase pair

Table I
REORDERING RULES EXTRACTED FROM FIGURE 1

| Num | Source | Target | POS tags |
|---|---|---|---|
| 1 | X1 的 X2 | X2 of X1 | NN NN NN NN |
| 2 | X1 后 X2 | X2 after X1 | NN NN VV NN |
| 3 | X1 后 确定 X2 的 X3 | decide X3 of X2 after X1 | NN NN NN NN NN NN |

we could calculate the frequency probabilities of it. Since the second step only matches the improved reordering rules whose number is small compared to the ordinary hierarchical rules, its time and space complexity is quite low.

## IV. USING IMPROVED REORDERING RULES IN DECODING

Our decoder is a re-implementation of Hiero [1][2]. The decoder is a CKY-style parser with beam search for every span in a chart. Given an input sentence $f$, the decoder searches for the most probable derivation which yields the translation $\hat{e}$ according to the following decision rule:

$$\hat{e}(f) = \underset{e \in D^*}{\operatorname{argmax}} \{ \sum_i \lambda_i \phi_i(f, e) \} \qquad (1)$$

where $D^*$ is the derivation set, $\phi_i(f, e)$ is the cost of one feature and $\lambda_i$ is the weight of one feature. Our decoder uses the same feature as Hiero[1].

Since improved reordering rules have the same form as the ordinary hierarchical rules, it is easy to incorporate them in hierarchical phrase-based translation systems. We add the improved reordering rules to the rule set and match them with no span length limit. In other words, we just treat them as additional hierarchical rules which could be used in every span.

[2]We have tested this decoder on various occasions to show that its performance is equivalent to Hiero.

When matching the improved reordering rules, we require that the POS tags of boundary words in the variable parts of the matching span should be the same as the variable boundary words' POS tags associated with the improved reordering rules. We refer to this as *POS constraint*. This constraint tries to prevent some potential wrong matching.

To speed up the rule matching, we use the dotted-rule matching in Earley-style parser[11] to get the matching hierarchical rules of a span, both for the ordinary rules and the improved reordering rules. This technique speeds up the matching process, especially for long spans.

Since the number of our reordering rule set is quite small and we use dotted-rule matching, incorporating the reordering rules will only add little cost. In our experiments, the average speed of the decoder drops from 16.4 words/sec to 13.9 words/sec.

## V. EXPERIMENTS

### A. Data

Our experiments were carried on Chinese-to-English translation in the news genre. We use FBIS dataset as the training corpus, NIST 2003 MT Evaluation (MT 03) as the development data, MT05 and MT08 as the test data. These files were segmented by ICTCLAS toolkit [12], and POS-tagged by a POS tagger implemented according to [13] and trained on Chinese Penn Treebank[7] . The bilingual corpus was word-aligned by Giza++ [14]. The ordinary

Table II
STATISTICS OF DEVELOPMENT AND TEST DATA(OOV STANDS FOR
UNKNOWN WORDS)

| Set | sentences | source words | OOV |
|------|-----------|--------------|------|
| MT03 | 919 | 26834 | 1864 |
| MT05 | 1082 | 33444 | 2058 |
| MT08 | 1357 | 33672 | 2453 |

Table IV
BLEU[%] SCORES IN EXPERIMENTS(**:SIGNIFICANCE AT THE 0.01
LEVEL)

| System | MT03 | MT05 | MT08 |
|--------|------|------|------|
| baseline | 28.66 | 28.02 | 19.75 |
| +reordering | 29.05** | 28.48** | 20.44** |

hierarchical rules and the improved reordering rules are both extracted from this corpus. We also use a 4-gram SRI language-model[15] trained on Xinhua portion of the Gigaword corpus. Detailed statistics of the development and test data are shown in Table II.

In the baseline system, we use only the ordinary hierarchical rules. The ordinary rules we extracted obey the same constraints as in [1]. We set the search depth to 10. In the contrast experiments, we add the improved reordering rules to the rule set.

The statistics of the ordinary rules (**ordinary**) and the improved reordering rules (**reordering**) are listed in Table III. Since some of the improved reordering rules may already exist in the ordinary rule set, we also show the number of unique rules (**unique**) of the improved reordering rules (Not found in the ordinary rule set). It accounts for $80.1\%$ of the whole improved reordering rule set. We also include the unique source side number(**source side**) to show how many reordering patterns there are.

We could see from the table that the number of the improved reordering rules is quite small compared with the number of the ordinary hierarchical rules. The results also show that most of the improved reordering rules are not extracted by the ordinary extraction process.

*B. Results*

To test the effect of the improved reordering rules, we run various kinds of experiments.

The first experiment (**+reordering**) is to compare the results of adding improved reordering rules with results of the baseline system. Note that the improved reordering rules are applied in every span during decoding. When matching reordering rules, we obey the POS constraint in Section IV.

The experiments' results are shown in Table IV. We use NIST-BLEU (case-insensitive) as the evaluation metric and methods in [16] for statistical significance tests. Results show that adding improved reordering rules improves BLEU scores in all NIST evaluation set and the results are statistically significant. Thus, we could see that overall, the

improved reordering rules improve the translation quality of the hierarchical phrase-based translation system.

The BLEU gains in three MT evaluation set are different. The main reason is that the sentence number of the three set are different (shown in Table II). Thus, the matching rule number are different. Our method is more helpful in long sentences which needs to be reordered non-locally.

We also test the effects of various kinds of rules on the MT03 set. In the first experiment, we eliminate the improved reordering rules extracted by replacing subrules (**-tree rule**). In the second experiment, we eliminate the reordering rules with consecutive variables(**-consecutive variables**). (Note that there are overlap rules in these two kinds.) The BLEU scores after making these changes in the reordering rule set are shown in the Table V. We could see that the rules with consecutive variables and the rules extracted by analyzing the tree structure both help reordering. We also found that in the translation results, these rules are applied to both local and non-local reordering.

We further test whether POS constraint is useful in decoding. The result is shown in Table V (**-POS constraint**). We could conclude that after applying the POS constraint, the reordering rules could capture more accurate reordering patterns and eliminate ambiguities.

Table VI shows an example sentence where the improved reordering rules capture the reordering pattern that ordinary rules fail to describe. (To save space, we only show part of the sentence.) Despite the existence of an out-of-vocabulary word, the correct reordering is made by adding improved reordering rules. It uses the following improved reordering rule which do not exist in the ordinary rule set.

- X → < 在 X1 X2 , X2 in X1 , NR NR NR NN >

However, the BLEU score gains are not as high as we expected. By analyzing the translation results, we found that although some long pieces of sentences are reordered correctly after adding the improved reordering rules, there are still errors in rule matching for long spans, which hurt

Table III
STATISTICS OF EXTRACTED RULES

| Rule Set | number | unique | source side |
|----------|--------|--------|-------------|
| ordinary | 31374474 | - | 21099125 |
| reordering | 105529 | 84591 | 72694 |

Table V
BLEU[%] SCORES(UNCASED) AFTER REMOVING DIFFERENT KINDS OF
REORDERING RULES AND CONSTRAINTS

| Condition | BLEU score |
|-----------|------------|
| baseline+reordering | 29.05 |
| -tree rule | 28.86 |
| -consecutive variables | 28.92 |
| -POS constraint | 28.35 |

Table VI
TRANSLATION COMPARISON BEFORE AND AFTER ADDING THE REORDERING RULES

| Source sentence | 以/MSP 保证/VV 在/P 科特迪瓦/NR 法国/NR 侨民/NN 的/DEC 安全/NN 。/PU |
|---|---|
| Baseline system | to guarantee security in 科特迪瓦 French nationals . |
| Adding reordering rules | to guarantee the safety of French nationals in 科特迪瓦 . |
| Reference Sentence | to ensure the safety of the French nationals in Cote d'Voire. |

the performance. This is probably due to the tagging errors, word alignment errors of long sentences, and the rules which contain too little lexical evidences. In our future work, beside POS constraints, we will try other methods to guide the rule matching and filter the ambiguous rules, such as phrase chunking.

## VI. CONCLUSIONS AND FUTURE WORK

In this paper, we propose an improved reordering rule extraction method and use these rules to improve the reordering of hierarchical phrase-based translation. These rules could be easily incorporated in a hierarchical phrase-based translation system with slight cost gain, and could be used in span of any length during decoding. Our experiments in Chinese-English translation show significant improvements over the baseline system for various NIST evaluation sets. In the future work, we plan to further improve the extraction process and investigate more feature functions to guide the incorporation of improved reordering rules to hierarchical phrase-based translation systems. We will also use these reordering rules in other phrase-based translation systems to help reordering.

## ACKNOWLEDGMENT

## REFERENCES

[1] D. Chiang, "Hierarchical phrase-based translation," *Computational Linguistics*, vol. 33, no. 2, pp. 201–228, 2007.

[2] C. Wang, M. Collins, and P. Koehn, "Chinese syntactic reordering for statistical machine translation," in *Proceedings of the 2007 Joint Conferences on Empirical Methods in Natural Language Processing and Computational Natural Language Learning(EMNLP-CoNLL)*, June 2007.

[3] P. Koehn, F. J. Och, and D. Marcu, "Statistical phrase-based translaiton," in *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*, May 2003.

[4] C. Tillman, "A unigram orientation model for statistical machine translation," in *Proceedings of HLT-NAACLs*, May 2004.

[5] B. Zhao and Y. Al-onaizan, "Generalizing local and non-local word-reordering patterns for syntax-based machine translation," in *Proceedings of the Conferences on Empirical Methods in Natural Language Processing*, Oct. 2008.

[6] F. J. Och and H. Ney, "The alignment template approach to statistical machine translation," vol. 30, no. 4, 2004, pp. 417–449.

[7] N. Xue, F. Xia, F.-D. Chiou, and M. Palmer, "The penn chinese treebank: phrase structure annotation of a large corpus," in *Natural Language Engineering*, vol. 11, 2005, pp. 207–238.

[8] M. Popovic and H. Ney, "Pos-based word reordering for statistical machine translation," in *Proceedings of the Fifth International Conference on Language Resources and Evaluation*, May 2006.

[9] J. M. Crego and N. Habash, "Using shallow syntax inforamtion to improve word alignment and reordering for smt," in *ACL 2008 Third Workshop on Statistical Machine Translation (WMT'08)*, June 2008, pp. 53–61.

[10] D. Xiong, Q. Liu, and S. Lin, "Maximum entropy based phrase reordering model for statistical machine translation," in *Proceedings of COLING-ACL 2006*, July 2006.

[11] J. Earley, "An efficient context-free parsing algorithm," in *Communications of the ACM*, vol. 13, no. 2, 1970, pp. 94–102.

[12] X.-Q. C. Hua-Ping Zhang, Qun Liu and H.-K. Yu, "Chinese lexical analysis using hierarchical hidden markov model," in *Proceedings of the second SIGHAN workshop on Chinese language processing*, vol. 17, 2003, pp. 63–70.

[13] M. Collins, "Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms," in *EMNLP*, July 2002, pp. 1–8.

[14] F. J. Och and H. Ney, "Improved statistical alignment models," in *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, Oct. 2000, pp. 440 –447.

[15] Stolcke and Andreas, "Srilm - an extensible language modeling toolkit," in *Proceedings of the International Conference on Spoken Language Processing*, Sept. 2002, pp. 311–318.

[16] M. Collins, P. Koehn, and I. Kucerova, "Clause restructuring for statistical machine translation," in *Proceedings of the 43th Annual Meeting of the Association for Computational Linguistics*, June 2005.