

統計翻訳における木構造制約の導入

山本 博史[†] 大熊 英男[†] 隅田英一郎[†]

[†] 〒 619-0288 京都府相楽郡精華町光台 2-2-2,
(独) 情報通信研究機構,
ATR 音声言語コミュニケーション研究所

E-mail: †{hirofumi.yamamoto,hideo.okuma,eiichiro.sumita}@nict.go.jp

あらまし 機械翻訳において、構文情報は非常に有用であり、近年ではこの情報を統計翻訳にも取り入れる試みがなされている。構文情報を統計翻訳に用いる場合、その情報をコーパスから統計的に学習することになるが、それを用いない場合と比較して大量のパラメータも学習する必要が生じる。大量のパラメータの学習を行うためには同じく大量の学習コーパスを必要とし、そのためにデータスパースネスの問題が生じる。このデータスパースネスの問題を避けるために、本稿では次の2点の仮定に基づく学習を必要としない構文情報モデルを提案する。翻訳原言語での関係のある単語は、翻訳先でも関係がる。関係を表すアークは交差しない。本モデルは SSMT2007 英中翻訳タスクのデータを用いた実験で、従来法より 1.9 ポイント高い BLEU 値 (31.3 から 33.2) と、4.9%低い WER (69.2%から 64.3%) を示し、有効性が確認できた。

キーワード 統計翻訳, 構文情報, データスパースネス

”Trainig-free” Tree Structure Model for SMT

Hirofumi YAMAMOTO[†], Hideo OKUMA[†], and Eiichiro SUMITA[†]

[†] National Institute of Information and Communications Technology,
ATR Spoken Language Translation Research Laboratories
2-2-2 Hikoridai, Seika-cho, Soraku-gun, Kyoto 619-0288 Japan.

E-mail: †{hirofumi.yamamoto,hideo.okuma,eiichiro.sumita}@nict.go.jp

Abstract For machine translation, syntax information is very useful. Recently, it is tried to introduce this information to statistical machine translation. In statistical machine translation, syntax information is extracted from corpus as statistical models through training. In this statistical models, larger number of parameters than conventional statistical translation models must be trained. Larger number of parameters requires larger training data. Request for larger training data makes data sparseness problem severer. To avoid this data sparseness problem, we proposed new syntax information introduction method. In proposed method, syntax information is modeled using following rules. The first is that word-to-word relation in translation source sentence is also kept in translation target sentence. The second is that word-to-word dependency arcs do not cross. This model is ”training-free” and has no data sparseness problem. In our experiments using SSMT2007 English-to-Chinese limited track data, proposed method result in 1.9 points improvements in BLEU (from 31.3 to 33.2), and 4.9% lower WER (from 69.2% to 64.3%) compared with base line conditions.

Key words Statistical Machine Translation, Syntax information, Data sparseness

1. はじめに

近年、機械翻訳として統計翻訳 (SMT), 特にフレーズベース統計翻訳 (PBSMT) [1] [2] [3] [4] が広く使われはじめ

ている。PBSMT における最も大きな問題点の一つとしてフレーズに並び替え (特に対極的な) がある。その理由は、PBSMT における並び替えモデルは単に並べ替えの際に何

単語先(後)に移動させるかの距離に依存したモデルであるためである。この問題を解決するために、構文情報を統計翻訳に導入する試みが数多くなされてきた。これらの試みは用いる翻訳原言語、翻訳先言語の構文情報の内どれを用いるかで大きく三つに分類される。一つ目は翻訳原言語、翻訳先言語双方の構文情報を用いるもので、tree-to-tree 翻訳 [5] [6] [7] [8] [9] [10] と呼ばれる。二つ目は翻訳先言語のみの構文情報を用いるもので、string-to-tree 翻訳 [11] [12] と呼ばれる。三つ目は tree-to-string 翻訳 [13] [14] と呼ばれ、翻訳原言語のみの構文情報を用いる。

これらの手法のように、構文情報を用いることで統計翻訳の性能を向上させることができる。しかしながら、構文情報をモデルに導入することによって学習すべきパラメータの数は増大する(特に tree-to-tree 翻訳では)。パラメータ数の増大は統計翻訳の訓練データであるパラレルコーパスに対し、質と量両面での要求が大きくなることを意味している。質の面ではそもそもデータスパースネスの問題が厳しい PBSMT に対し、さらに拍車をかけることになる。また、質の面から考えると学習データたるパラレルコーパスの対訳文には原言語の構造が反映されている必要があるが、同一の意味でも何種類かの異なった構文に翻訳が可能である。これは、構造の反映のさせ方に必ずしも一貫性がないことを意味しており、これもまたデータスパースネスの問題にさらに拍車をかけることになる。

このように、統計翻訳に対する構文情報の導入はモデルの表現能力を高めるものの、モデルパラメータの学習の面では深刻なデータスパースネスの問題を引き起こす。そこで、我々はこのデータスパースネスの問題を回避するために、パラメータ学習の不要な構文情報モデルである木構造制約モデルの導入を試みる。

まず第二章では最も単純なケースとして、一対一の単語対応が取れる条件での木構造制約モデルの説明を行い、第三章でそれをフレーズベースの翻訳モデルに拡張する。第四章ではこのモデルを用いてデコーディングを行う際のアルゴリズムを示す。第五章では提案モデルに対する実験条件とその結果について述べ、第六章でまとめを行う。

2. 訓練データを用いない木構造制約モデル

まず、最も単純なケースとして、翻訳原言語文、翻訳先文の全ての単語が 1 対 1 対応をしている場合を考える。翻訳原言語文中の単語 s_i は翻訳先では S_i に翻訳されるものとした場合、翻訳原言語文 s_1, s_2, \dots, s_N は単語セット $\{S_1, S_2, \dots, S_N\}$ の語順を並び替えたものとして翻訳される。この場合翻訳先文の可能な組み合わせ数は $N!$ となる。木構造制約モデルの目的は、この $N!$ 通りの探索空間を縮めるような制約を与えることにある。

木構造制約モデルは次に示す二つのルールに従いながら翻訳が可能であるという仮定に基づいており、パラメータの訓練を必要としない。

- ルール 1: 翻訳原言語単語 s_i が s_j と依存関係等の関係を持つならば、翻訳先言語単語 S_i もまた S_j と関係を持つ。

- ルール 2: 単語間の関係を表すアークは交差しない。

上記のルールを満たしているかどうかのチェックには、bracketed sentence を用いる。英文 "This is a pen." をパースした結果として構文木で (S1 (S (NP (DT This)) (VP (AUX is) (NP (DT a) (NN pen)))) (. .)) が得られる。ここから全ての句ラベルを取り除くことによって bracketed sentence ((This) ((is) ((a) (pen)))) (.) が得られる。この制約ルールを適用した場合、翻訳原言語の bracketed sentence ((a b) (c d)) から得られる翻訳先言語文は、次のどれかに限られることになる。[A B C D], [A B D C], [B A C D], [B A D C], [C D A B], [C D B A], [D C A B], [D C B A]。ここで、翻訳先言語単語 A, B, C, D はそれぞれ翻訳原言語単語 a, b, c, d の対訳語である。翻訳先言語文として [A C B D] を考えてみる。翻訳原言語の bracketed sentence から、a と b, c と d は関係を持つ。これに対し、ルール 1 を適用すると A と B, C と D もまた関係を持つ。しかしながら、この関係を単語列 [A C B D] にあてはめるとそれぞれの関係を示すアークが交差することになり、ルール 2 を満たすことができない。本制約を適用した場合、翻訳先文は次のようにして得られることになす。翻訳原言語の bracketed sentence を表す木構造の各ノードに対し、その直下のサブツリー(またはリーフである単語)どうしの順序を入れ替える。たとえば、図 1 において、ノード 1 に対してのみ入れ替えを行うことによって、翻訳先言語文 [B A C D] が得られ、ノード 2 と 3 に対して入れ替えを行うことによって [D C A B] が得られる。

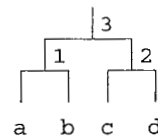


図 1 bracketed sentence の例

実際の文の例としては、英文 (He (eats (large bread) quickly) .) に対する日本語対訳、(彼は ((大きなパンを) 早く 食べる) .) があげられ、SVO と SOV と語順に大きな違いがあるにもかかわらず、上記のルールを満たしながら翻訳ができています。

木構造制約モデルを用いない場合 ((a b) (c d)) に対する可能な翻訳先文の組み合わせは $4! = 24$ である。一方、本モデルを導入した場合は 8 に減少している。N 単語からなるバイナリ bracketed tree の場合、本モデル導入時の組み合わせ数は 2^{N-1} となる。この理由はこのバイナリ木のノードの数は $N-1$ であり、それぞれのノードに対し、入れ替えを行う、行わないの二つの選択枝があるためである。この組み合わせ数は、本モデルを導入しない場合の $N!$ に比べ、

非常に小さいものとなっており、実際 $N = 10$ の場合で約 $1 / 7,000$, $N = 20$ の場合で $1 / 2 \times 10^{12}$ である。より一般的に、バイナリ木でない場合の本モデル導入時の組み合わせ数は $\prod_{i=1}^n (B_i!)$ である (ここで n は木に含まれるノードの数、 B_i は i 番目のノードの枝の数を表す)。

式 (1) は、本モデルを用いない場合の統計翻訳を表す式であり、 $P(f|e)$, $P(e)$ はそれぞれ翻訳モデル、言語モデルを表している。

$$\operatorname{argmax}_e P(e|f) = \operatorname{argmax}_e P(e)P(f|e) \quad (1)$$

これに対し、提案モデルを用いた場合は新たな項 $P(e|T)$ が追加され、次のような式で表されることになる。

$$\operatorname{argmax}_e P(e)P(f|e)P(e|T) \quad (2)$$

ここで $P(e|T)$ が本木構造制約モデルであり、 T は翻訳原言語文の bracketed tree を表している。 $P(e|T)$ の値は e がモデルの制約を満たす時は 1 であり、そうでなければ 0 である。

3. フレーズベースモデルへの拡張

前章では 1 対 1 単語対応の場合の木構造制約モデルについての説明を行った。本章ではこれをフレーズベースモデルに対して適用できるように拡張する。通常、単語アライメントは n 対 m (0 対 m , n 対 0 を含む) である。しかしながら、フレーズベースモデルでは、フレーズ対フレーズのアライメントはたとえそれぞれのフレーズに含まれる単語数が異なっていたとしても、常に 1 対 1 対応となる。このため、前章の 1 対 1 単語対応のルールをおおむねそのままフレーズ対フレーズの対応に当てはめることができる。フレーズ ph_i が単語 s_n を含み、フレーズ ph_j が単語 s_m を含むものとする。ここで、単語 s_n と単語 s_m が関係を持つならば、フレーズ ph_i とフレーズ ph_j も関係を持つと定義する。これにより、前章のルールは次のようにフレーズに拡張することができる。

- ルール 1: 翻訳原言語フレーズ ph_i が ph_j と依存関係等の関係を持つならば、翻訳先言語フレーズ PH_i もまた PH_j と関係を持つ。
- ルール 2: フレーズ間の関係を表すアークは交差しない。

ここで PH_n は翻訳原言語フレーズ ph_n の対訳フレーズを表すものとする。

このフレーズに拡張したルールを適用した場合、1 対 1 単語対応の場合と同様、翻訳原言語の bracketed sentence を表す木構造の各ノードに対し、その直下のサブツリー (またはリーフである単語) どの順序を入れ替えることによって翻訳先言語文が得られる。ただし、入れ替えを行えるノードには制限があり、フレーズの一部のみを含むノードは入れ替えを行うことができない。たとえば図 2 のように、bracketed tree $((abc)((de)(fg)))$ で、 bcd

がフレーズ ph を構成している場合を考える。この場合ノード 1 はフレーズ ph の一部である bc を含むため入れ替えができない。同様に、ノード 2、4 も入れ替えができない。一方ノード 3 はフレーズを含まず、ノード 5 はフレーズ全体を含んでいるため入れ替え可能である。たとえばノード 2 に対し入れ替えを行った場合、フレーズ ph は翻訳先言語では二箇所に分割されることになり、フレーズベースモデルにおけるフレーズ対フレーズの対応が 1 対 1 であることに反することになる。結果として、この bracketed tree の対訳としては $[APH EFG]$, $[APHEGF]$, $[GF EPHA]$, $[FGE PHA]$ のみが許されることになる。ここで、 PH は ph の対訳フレーズを表すものとする。

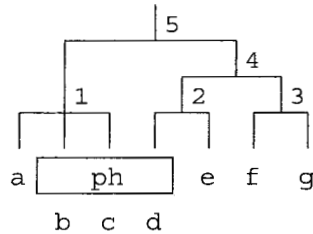


図 2 フレーズを含む木構造の例

4. 木構造モデルのデコーディングアルゴリズム

本章では木構造制約モデルを target side left-to-right デコーダに導入する場合のアルゴリズムについて記述する。このデコーダでは翻訳先言語文は左から右へ (文頭から文末に向かって) 順に生成されていく。翻訳先言語文仮説を右に伸ばすためのフレーズが新たに選択された場合、そのフレーズに対応する翻訳原言語の単語に対してビットが立てられる。そして、翻訳原言語の全ての単語に対するビットが立てられた時に翻訳先言語文仮説は文末に達したと判断される。木構造制約モデルをこの流れに組み込むためには翻訳先言語文フレーズが新たに選択されるたびに木構造制約モデルの制約を満たしているかどうかのチェックをする必要がある。

このチェックアルゴリズムの説明の前に翻訳原言語の bracketed tree のサブツリーを NOT-TRANSLATED, TRANSLATED, TRANSLATING, NG の 4 種類に分類しておく。

- もし、サブツリーがリーフである単語のみからなっており、かつ全ての単語が未翻訳 (ビットが立っていない) ならば、そのサブツリーは NOT-TRANSLATED である。
- もし、サブツリーが NOT-TRANSLATED サブツリーのみからなっているならば、そのサブツリーもまた NOT-TRANSLATED である。
- もし、サブツリーがリーフである単語のみからなっており、かつ全ての単語が翻訳済み (ビットが立っている) ならば、そのサブツリーは TRANSLATED である。

● もし、サブツリーが **TRANSLATED** サブツリーのみからなっているならば、そのサブツリーもまた **TRANSLATED** である。

● もし、サブツリーがリーフである単語のみからなっており、未翻訳の単語、翻訳済みの単語双方を含むならば、そのサブツリーは **TRANSLATING** である。

● もし、サブツリーが **TRANSLATED**, **NOT-TRANSLATED** 双方のサブツリーを含むならば、そのサブツリーは **TRANSLATING** である。

● もし、サブツリーが一つだけ **TRANSLATING** サブツリーを含むのならば、そのサブツリーは **TRANSLATING** である。

● もし、サブツリーが二つ以上の **TRANSLATING** サブツリーを含むのならば、そのサブツリーは **NG** である。

● もし、サブツリーが **NG** サブツリーを含んでいるならば、そのサブツリーもまた **NG** である。

デコーディング中に **NG** サブツリーが生成されたならば、その翻訳仮説は木構造制約モデルを満たすことができない。翻訳原言語側のサブツリーが単語列 $[x_1, x_2, \dots, x_n]$ からなるとする。この時、このサブツリーの対訳はそれの対訳語のセット $\{X_1, X_2, \dots, X_N\}$ の語順を並べ替えたものとして得られる。もし、途中で他の単語が割り込んだ場合、ルール 2 が満たせなくなるためである。このため、**TRANSLATING** サブツリーの最後の翻訳済み単語に続く単語は、このサブツリーの中の未翻訳の単語でなければならない。これは次に翻訳される単語は **TRANSLATING** サブツリーの中の未翻訳の単語から選ばなければならないことを意味している。たとえば、翻訳済単語 a, b と、未翻訳単語 c, d, e, f からなるサブツリー $((ab)((cd)(ef)))$ が翻訳文仮説に含まれるならば、次に翻訳されるべき単語は c, d, e, f のうちのどれかでなければならない。従って、**TRANSLATING** サブツリーが二つ以上含まれるならば、この条件を満たすことができなくなる。

翻訳原言語フレーズ ph の対訳フレーズ PH を翻訳文仮説 (生成中) に後続させる場合を考える。もとの翻訳文仮説を (S_1, S_2, \dots, S_i) とし、これは木構造制約モデルを満たしているものとする。これにフレーズ PH を後続させた仮説が **TRANSLATING** を二つ以上生成させない (すなわち木構造制約モデルを満たす) ためには、次の条件のうちのどちらかを満たす必要がある。ここで、 T は (S_1, S_2, \dots, S_i) に含まれる最小の **TRANSLATING** サブツリーを表すものとする。

(1) フレーズ PH を後続させた後でも T が **TRANSLATING** であり、かつ別の **TRANSLATING** サブツリーが生成されていないこと。

(2) フレーズ PH を後続させることによって、 T が **TRANSLATED** サブツリーとなること。

フレーズ ph が T の未翻訳の部分に含まれていることは

条件 1 に対して必要十分である。また、フレーズ ph が T の未翻訳の部分を含んでいることは条件 2 に対して必要十分である。このことより、新たに生成される翻訳仮説が木構造制約モデルを満たしているかのチェックは次の手順で行うことができる。

(1) 生成されている翻訳仮説に対し、それに含まれる最小の **TRANSLATING** サブツリーを記憶しておく。

(2) 旧翻訳仮説に対し、新たなフレーズ PH を後続させる場合、その対訳である翻訳原言語フレーズ ph と、旧翻訳仮説の記憶されている最小の **TRANSLATING** サブツリーの未翻訳部分との比較を行う。もし、 ph が未翻訳部分に含まれる、 ph が未翻訳部分を含むならば、旧翻訳仮説に PH を後続させたあらたな仮説を生成し、それに含まれる最小の **TRANSLATING** サブツリーを更新する。そうでなければ、この仮説は棄却される。

5. 実 験

5.1 評価尺度

提案法の評価のためには四つの評価尺度 WER, PER, BLEU [15], NIST [16] を用いた。評価を行う前に各評価尺度に対する提案法の有効性を推測してみる。

● WER: この尺度は大域的な単語順序の入れ替えを考慮することができる。そのため提案法はこの尺度に対して有効に働くと予想される。

● PER: この尺度は基本的には語順を考慮することができない。従ってこの尺度に対して提案法は有効ではないと予想される。

● BLEU: この尺度は ngram に着目するため、中距離の単語順序の入れ替えを考慮することができる。たとえばレファレンス翻訳 translation (w_1, w_2, \dots, w_n) に対し、翻訳結果が $(w_1, w_2, \dots, w_{j-1}, X, w_{j+1}, \dots, w_n)$ である場合 WER, BLEU は共に高い値を示す。しかしながら、翻訳結果 $(w_{j+1}, \dots, w_n, X, w_1, w_2, \dots, w_{j-1})$ に対し BLEU は同じく高い値を示すのに対し、WER の値は 0 となる。したがって提案法は BLEU に対し有効ではあるが、WER ほどではないと予想される。

● NIST: この尺度も BLEU 同様 ngram に着目する。しかしながら、低次の ngram に対する重みが BLEU よりおおいいため提案法の有効性は BLEU よりも低いと予想される。

5.2 英日ニュース翻訳実験

まず提案法の性能評価のために英日ニュース翻訳実験を行った。提案法を用いるためにはまず翻訳原言語文の bracketed sentence が必要となり、このために翻訳原言語文をパズルする必要がある。この時パーズングエラーによる性能劣化が予想される。本実験では提案法の性能評価を行うとともに、パーズングエラーによる性能劣化の評価を行うことも目的とする。パーズングエラーによる性能劣化の評価のために、提案法に対しては自動でパーズングを行った

結果と正しい(人手であたえた) パーズ結果を用いた場合の二通りに対する評価を行った。

実験コーパスとしては読売新聞、ロイター [17] およびウォールストリートジャーナルを訓練コーパスとして用いた。それぞれのデータサイズは 145K, 57K, 14K である。またウォールストリートジャーナルから 1,787 文をデベロップメントセットとして、同じく 1,787 を評価セットとして用いた。実験に用いたウォールストリートジャーナル文はベンツリーバンクコーパス [18] [19] に含まれているものであり、人手によるパーズツリーが与えられている。これらのコーパスの詳細については、表 1 に示されている。

表 1 英日ニュース翻訳実験コーパス

	# of sent.	Total words	# of entry
E/J Train	216K	5.6M/7.0M	82K/61K
E/J Dev	1,787	43K/43K	7,836/7,229
E Eval	1,787	64K	7,316

フレーズベース翻訳モデルの訓練には GIZA++ [20] を、言語モデルの訓練には SRI language model tool kit [21] を用いた。言語モデルは単語トライグラムで、Kneser-Ney ディスカウンティング [22] で平滑化を行った。デコーディングパラメータの最適化には minimum error training [23] を用い、BLEU に対して最適化を行っている。また、翻訳原言語の bracketed tree の抽出には Charniak パーザー [24] を用いた。

実験では三つの条件での比較を行った。Base-line は提案法を用いない場合、Chariniak は bracketed tree の抽出に Charniak パーザーを用いた場合、Oracle はベンツリーバンクの木構造を用いた場合である。デコーダは我々が独自に開発した Pharaoh [25] 互換デコーダ CleopATRA を木構造制約モデル用に改造したものをを用いた。この際のパラメータは全て共通で、Base-line の条件で最適化を行ったものをを用いた。表 2 に各条件での評価結果を示す。Chariniak パーザーを用いた条件 Chariniak で WER は約 4% の改善、BLEU では約 0.6 の改善であった。各評価基準に対する改善幅は前節での予想と一致しており、WER が一番で BLEU がそれに続く。正解木構造を用いた Oracle の結果と Chariniak では大きな違いはなく、提案法に対しては Charniak パーザーの精度は十分であるといえる。この時評価文セットに対する bracketed tree は Chariniak と Oracle で 60% が同じであった。

表 2 英日ニュース翻訳実験評価結果

	BLEU	NIST	WER	PER
Base-line	15.4	5.13	86.9	55.1
Chariniak	16.0	5.20	83.0	54.9
Oracle	16.1	5.20	82.8	54.8

5.3 英日特許翻訳実験

続いて先の実験とは異なるドメインである特許に対する翻訳実験を行った。特許翻訳実験コーパスに関する詳細情報を表 3 に示す。モデルの訓練、パラメータの最適化方法、デコーディングに関しては、ニュース翻訳実験で用いた方法と同じである。

表 3 英日特許翻訳実験コーパス

	# of sent.	Total words	# of entry
E/J Train	10G	273M/257M	797K/282K
E/J Dev	999	39K/37K	4,971/4,614
E Eval	556	20K	2,462

表 4 に実験結果を示す。ニュース翻訳実験の場合と同様 WER に対する改善が最も大きく 4.9% で、BLEU がそれについて 1.5 である。この実験結果から提案手法は異なるドメインに対しても有効であることが確認できた。

表 4 英日特許翻訳実験評価結果

	BLEU	NIST	WER	PER
Base-line	25.3	6.68	84.4	44.8
Proposed	26.8	6.73	79.5	44.8

5.4 英中翻訳実験

最後の異なる言語ペアとして英中翻訳実験を行った。実験に用いたコーパスは SSMT2007 [26] 英中リミテッドトラックで用いられたもので、その詳細を表 5 に示す。モデルトレーニング等の条件は英日実験の場合と同様であるが、パラメータの最適化のみに対しては評価セットをそのまま用いており、パラメータに関してクロズドの条件となっている。

表 5 SSMT2007 英中翻訳実験コーパス

	# of sent.	Total words	# of entry
E/C Train	835K	10.1M/9.6M	149K/94K
E Eval	505	11K	2,972

表 6 に実験結果を示す。なお、本実験での中国語レファレンスの数は 4(日本語リファレンスは 1) となっている。また評価の単位は文字(漢字)である。英日実験の場合と同様 WER に対する改善が最も大きく 4.9% で、BLEU がそれについて 1.9 である。この実験結果から提案手法は異なる言語ペアに対しても有効であることが確認できた。

表 6 SSMT2007 英中翻訳実験評価結果

	BLEU	NIST	WER	PER
Base-line	31.3	7.49	69.2	41.9
Proposed	33.2	7.61	64.3	42.0

6. ま と め

本稿では、フレーズベース統計翻訳 PBSMT に対して構文情報を導入するための手法を提案した。提案手法である木構造制約モデルは言語非依存の二つのルールに基づいている。一つ目は翻訳原言語における単語間の依存関係等は翻訳先言語でも保持されるというもので、二つ目はその単語間の関係を表すアークは交差しないというものである。木構造制約モデルの最大の特徴は訓練が不要であることであり、そのため従来の構文情報を用いたモデルで問題となっていたデータスパースネスの問題を生じることがない。本モデルは left-to-right デコーダに直接組み込むことが可能で、その際には単語の並び替えに関するあらたな制約として働く。本モデルは英中翻訳実験において BLEU で 1.9, WER で 4.9% の改善を示し、単語の大域的な並び替えに関して有効に働くことが WER における 4.9% という性能向上で確認することができた。

7. 謝 辞

本研究の遂行にあたりまして、実験データとして用いました特許コーパスの提供および使用許諾をくださいました日本特許情報機構 (JAPIO) 様に感謝をいたします。

文 献

- [1] Daniel Marcu, William Wong, "A phrase-based, joint probability model for statistical machine translation," Proc. EMNLP-2002, pp.133-139, 2002.
- [2] R. Zens, F. J. Och, H. Ney, "Phrase-based statistical machine translation," 25th German Conference on Artificial Intelligence, Sep 2002.
- [3] P. Koehn, F. J. Och, D. Marcu, "Statistical phrase-based translation," Proc. HLT-NAACL, pp. 127-133, 2003.
- [4] F. J. Och, H. Ney, "The alignment template approach to statistical machine translation," Computational Linguistics, 30(4), pp417-449, 2004.
- [5] Dekai Wu, "Stochastic inversion transduction grammars and bilingual parsing of parallel corpora," Computational Linguistics, 23(3), pp.377-403, 1997.
- [6] H. Alshawi, S. Bangalore, S. Douglas, "Learning dependency translation models as collections of finite-state head transducers," Computational Linguistics, 26(1), pp.45-60, 2000.
- [7] D. Melamed, "Statistical machine translation by parsing," Proc. ACL, pp. 653-660, 2004.
- [8] D. Chiang, "A Hierarchical Phrase-Based Model for Statistical Machine Translation," Proc. ACL2005, pp. 263-270, 2005.
- [9] Y. Ding, M. Palmer, "Machine translation using probabilistic synchronous dependency insert grammars," Proc. ACL, pp. 541-548, 2005.
- [10] 今村 賢治, 大熊 英男, 隅田 英一郎, "句に基づく構文トランスファ方式統計翻訳," 情報処理学会論文誌, Vol.48, No. 4, Apr, 2007.
- [11] K. Yamada, K. Knight, "A syntax-based statistical translation model," Proc. ACL, pp. 523-530, 2000.
- [12] Daniel Marcu, Wei Wang, Abdessamad Echihabi, Kevin Knight, "SPMT: Statistical Machine Translation with Syntactified Target Language Phrases," Proc. EMNLP-2006, pp. 44-52, 2006.
- [13] C. Quirk, A. Menezes, C. Cerry, "Dependency treelet translation: Syntactically informed phrasal SMT," Proc. ACL, pp. 271-279, 2005.
- [14] Yang Liu, Qun Liu, Shouxun Lin, "Tree-to-String Alignment Template for Statistical Machine Translation," Proc. ACL2006, pp. 609-616, 2006.
- [15] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," Proc. ACL, 2002.
- [16] G. Doddington, "Automatic evaluation of machine translation quality using n-gram co-occurrence statistics," Proc. ARPA Workshop on Human Language Technology, 2002.
- [17] Masao Utiyama and Hitoshi Isahara, "Reliable Measures for Aligning Japanese-English News Articles and Sentences", ACL-2003, pp. 72-79, 2003 .
- [18] M. P. Marcus, B. Santorini, M. A. Marcinkiewicz, "Building a Large Annotated Corpus of English: The Penn Treebank," Computational Linguistics, 19(2):313-330.
- [19] Kiyotaka Uchimoto, Yujie Zhang, Kiyoshi Sudo, Masaki Murata, Satoshi Sekine and Hitoshi Isahara, "Multilingual Aligned Parallel Treebank Corpus Reflecting Contextual Information and Its Applications," Proc. MLR2004: PostCOLING Workshop on Multilingual Linguistic Resources, pp. 63-70, 2004.
- [20] F. J. Och, H. Ney, "A Systematic Comparison of Various Statistical Alignment Models," Computational Linguistics, No. 1, Vol. 29, pp. 19-51, 2003.
- [21] A. Stolcke, "SRILM - An Extensible Language Model Toolkit," Proc. ICSLP'02, 2002.
<http://www.speech.sri.com/projects/srilm/>
- [22] R. Kneser, H. Ney, "Improved backing-off for m-gram language model," Proceedings of the IEEE International Conference of Acoustic, Speech, and Signal processing. Vol. 1, pp. 181-184, 1995.
- [23] F. J. Och, "Minimum error rate training for statistical machine translation," Proc. ACL, 2003.
- [24] E. Charniak, "A Maximum-Entropy-Inspired Parser," Proc. NAACL-2000, pp.132-139, 2000.
- [25] P. Koehn, "PHARAOH: A beam search decoder for phrase-based statistical machine translation models," Proc. AMTA, 2004.
<http://www.isi.edu/publications/licensed-sw/pharaoh/SSMT2007/>
- [26] The Third Symposium on Statistical Machine Translation
http://mitlab.hit.edu.cn/EvaluationGuidelines_En.html