

# 統計的機械翻訳のための単語列ラベリングに基づく語順推定モデル

Distortion Model Based on Word Sequence Labeling for Statistical Machine Translation

日本放送協会 放送技術研究所ヒューマンインターフェース研究部専任研究員 **後藤 功雄**

## PROFILE

1997年NHK入局。2008年より情報通信研究機構に出向。2013年NHKに復帰。2014年京都大学大学院博士課程修了。博士（情報学）。自然言語処理の研究に従事。

## 1 はじめに

筆者らは統計的機械翻訳（SMT）において語順推定を改良する研究を情報通信研究機構にて実施した<sup>[1, 2]</sup>。本稿では、この研究成果について紹介する。SMTでは、訳語選択と語順推定という2つの処理が必要である。特に日英など語順が大きく異なる言語間では、語順の推定が難しい課題である。この課題を解決するために、これまでに語彙化語順推定モデル<sup>[3]</sup>、並べ替え制約<sup>[4]</sup>、プレオーダーリング<sup>[5]</sup>、階層フレーズベースSMT<sup>[6]</sup>、構文ベースSMT<sup>[7]</sup>などが提案されている。一般的に、原言語の構文構造は長距離の語順推定に有用である。しかし、多くの言語では高性能な構文解析器が利用できない。そこで、構文解析器を必要としない翻訳手法も必要である。フレーズベースSMT<sup>[8]</sup>は、構文解析器を必要としないSMT手法で広く利用されているものの1つである。我々は、フレーズベースSMTで利用する新しい語順推定モデルを提案する。

## 2 フレーズベースSMTでの語順推定モデル

フレーズベースSMT<sup>[8]</sup>では、目的言語文を文頭から文末へ連続的に生成していく。そのため語順推定モデルの役割は、翻訳処理中において最後に生成した目的言語の表現の右隣に生成すべき目的言語表現に対応する原言語文中の表現の位置を推定することである。例を図1に示す。図1では、「彼は」のみが翻訳されている状態

を表している。次に生成されるべき語は“bought”であるので、次に翻訳されるべきである原言語文中の語は「買った」である。そのため、語順推定モデルは、「買った」を含むフレーズを次に翻訳すべきフレーズの位置として推定する必要がある。



図1 日英翻訳の語順推定の例。四角の枠はフレーズを表す。

ここで、語順推定の課題をより正確に説明するために current position (CP) および subsequent position (SP) という2つの用語を導入する。CPは、生成中の目的言語文中でアラインメントされている右端の語に対応する原言語文中の位置である。SPは、次に生成する目的言語のフレーズ中でアラインメントされた左端の語に対応する原言語文中の位置である。フレーズベースSMTによる翻訳処理での語順推定モデルの役割は、原言語文中において最後に翻訳した位置（CP）が分かっている状態で次に翻訳する位置（SP）<sup>1</sup>を推定することである<sup>2</sup>。

SPの推定は難しい課題である。図2にCPとSPの例を示す。上付き数字は原言語の単語位置を示す。

図2において、(a)ではSPは8である。しかし、(b)ではCPの語は同じであるがSPは異なる。これらの例

- 1 正解は複数あり得るため、SPは1カ所だけとは限らない。
- 2 この定義は既存手法<sup>[8]</sup>の定義とは少し違いがある。既存手法では、フレーズ内の単語対応を考慮せずに、CPは最後に翻訳したフレーズの右端位置、SPは次に翻訳するフレーズの左端位置としている。

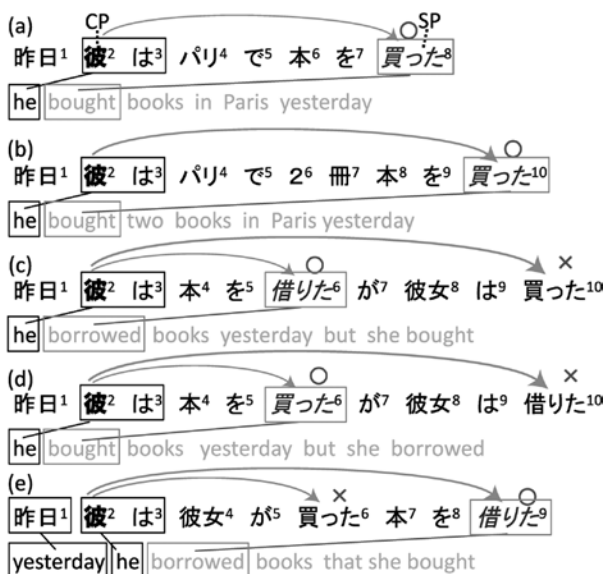


図2 日英翻訳でのCPとSPの例。CPはボールド、SPはイタリックで表されている。×はSP以外のSP候補を表している。

から距離はSPを決める本質的な要因ではないことが分かる。また、SPを推定するのにCPの語とその文脈だけでは十分ではないことが分かる。そのため、CPの語だけでなくSP候補の語も同時に考慮する必要がある。

図2において、(c)と(d)ではCPの語は同じであり、どちらにもSP候補として「借りた」と「買った」がある。(c)では「借りた」がSPであるが、(d)では「買った」がSPである。この違いの原因の1つはSP候補間の相対的な語順の違いである。そのため、SP候補の語だけでなく、SP候補間の相対的な語順も重要である。

図2において、(d)と(e)ではCPの語は同じであり、SP候補の「買った」と「借りた」の語順も同じである。しかし、(d)と(e)ではSPの語が異なる。この要因はCPやSP候補の語の文脈の違いである。そのため、近いSP候補を選択すれば良いわけではなく、CPやSP候補の文脈を考慮する必要がある。

すなわち、SPを推定するためには、CPの語、SP候補の語、SP候補間の相対的な語順、CPおよびSP候補の文脈を考慮する必要がある。

フレーズベースSMTの語順推定モデルはこれまでに、MSD語彙化語順推定モデル<sup>[3, 9, 10]</sup>や、CPもしくはSP候補の文脈を考慮する語順推定モデル<sup>[11, 12, 13]</sup>が提案されている。これらのモデルはいずれもSPをCPからの相対位置(距離と方向)で特定し、相対位置の確率を計算する。相対位置はいくつかに分類されて、

各分類(クラス)の確率を計算する。この場合、同じクラスに含まれる位置が同じ確率になってしまうという問題がある。この問題を軽減するために分類を細かくすると、各クラスに該当する学習データの量が少なくなると多くのクラスでフレーズや語の頻度が0になり、統計量の信頼性が低下してしまうという問題が起きてしまう。

### 3 提案手法

我々は、既存手法のように相対位置の確率を計算するというアプローチではなく、各SP候補を特徴付ける特徴量に基づいて各SP候補がSPである確率を直接計算するというアプローチを提案する。これによって、前記の問題を起こさずに各位置の確率を計算することができる。SP候補*j*がSPである確率は、

$$P(X=j|i, S) \quad (1)$$

と表すことができる。ここで、*i*はCP、*S*は原言語文、*X*はSPの確率変数である。

次に(1)式を実現する2つのモデルを説明する。

#### 3.1 ペアモデル

ペアモデルは、(1)式を実現する手法で、メインの提案手法のベースとなるモデルである。我々はこのモデルを最大エントロピー法で構築する。*s*を原言語の単語、 $s_1^i = s_1 s_2 \dots s_n$ を原言語文とする。文頭にBOS、文末にEOSの記号を追加する。すると原言語文*S*は $s_0^{n+1}$  ( $s_0 = \text{BOS}$ ,  $s_{n+1} = \text{EOS}$ )で表される。ペアモデルは次式により確率を計算する。

$$P(X=j|i, S) = \frac{1}{Z_i} \exp(\mathbf{w}^T \mathbf{f}(i, j, S, o, d)) \quad (2)$$

ここで、

$$o = \begin{cases} 0 & (i < j) \\ 1 & (i > j) \end{cases}, \quad d = \begin{cases} 0 & (|j-i|=1) \\ 1 & (2 \leq |j-i| \leq 5) \\ 2 & (6 \leq |j-i|) \end{cases}$$

$$Z_i = \sum_{j \in \{1 \leq j \leq n+1 \wedge j \neq i\}} \exp(\mathbf{w}^T \mathbf{f}(i, j, S, o, d))$$

であり、 $\mathbf{f}(\cdot)$ はバイナリ素性関数を要素とするベクトルで、 $\mathbf{w}$ はそれに対応する重みパラメータのベクトルである。 $Z_i$ は正規化項、*o*は*i*から*j*への方向、*d*は距

表 1 素性テンプレート

テンプレート
$\langle o \rangle, \langle o, s_{i+p} \rangle^1, \langle o, s_{j+q} \rangle^1, \langle o, t_i \rangle, \langle o, t_j \rangle, \langle o, d \rangle,$
$\langle o, s_{i+p}, s_{j+q} \rangle^2, \langle o, t_i, t_j \rangle, \langle t, t_{i-1}, t_i, t_j \rangle, \langle o, t_i, t_{i+1}, t_j \rangle,$
$\langle o, t_i, t_{j-1}, t_j \rangle, \langle o, t_i, t_j, t_{j+1} \rangle, \langle o, s_i, t_i, t_j \rangle, \langle o, s_j, t_i, t_j \rangle$

<sup>1</sup>  $p \in \{0\} \cup \{-2 \leq p \leq 2\}$

<sup>2</sup>  $(p, q) \in \{(p, q) \mid -2 \leq p \leq 2 \wedge -2 \leq q \leq 2 \wedge (|p| \leq 1 \vee |q| \leq 1)\}$

離クラスである。バイナリ素性関数は素性に一致した場合に1、それ以外で0を返す。表1に素性を構築する際に用いる素性テンプレートを示す。 $t_i$ は $s_i$ の品詞を表す。

(2)式で $i, j, S$ は素性関数で利用されている。これによって、CPの語、SP候補の語、CPおよびSPの文脈を同時に考慮して確率を計算することができる。

### 3.2 系列モデル

ペアモデルではSP候補間の相対的な語順を考慮できないという問題がある。この問題を解決するモデルを提案する。このモデルがメインの提案手法であり、系列モデルと呼ぶ。

図2の(c)と(d)で、「借りた」と「買った」はどちらも原言語文に出現している。ペアモデルは距離の違いによる影響を距離クラスである $d$ のみで考慮している。そのため、もしこれらの語が同じ距離クラスに出現した場合は、これらの語のCPからの距離の違いを扱うことができない。このようなケースは、学習時では適切に分離できない学習データとなり、翻訳時では適切に推定することが困難となる。

系列モデルは相対的な語順を考慮することで、この問題を解決する。これは、CPから各SP候補(SPC)の範囲を表すラベル系列を識別することで行う。各ラベル系列は1つのSPCに対応しているので、SPに対応するラベル系列を特定することができれば、SPを特定することができる。ラベル系列はC、I、Sという3種類のラベルを用いてCPからSPCの範囲を表す。ラベルCはCPを表し、ラベルSはSPCを表し、ラベルIはCPとSPCの間の位置を表す。図2(c)に対するラベル系列の例を図3に示す。SPCは各ラベル系列のIDとして用いる。

ラベル系列は相対的な語順を扱うことができる。例えば、図3のラベル系列IDが10のラベル系列は、「買った<sup>10</sup>」のSPCよりCPに近い位置に「借りた」が存在

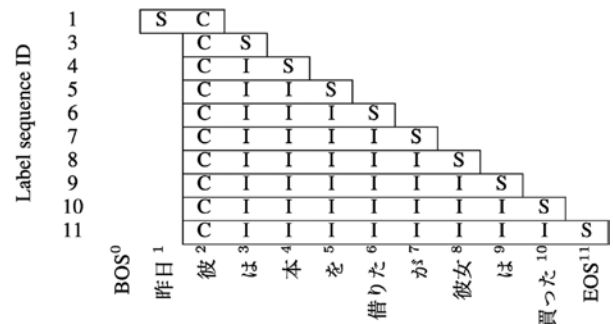


図3 ラベル系列の例。図2(c)の場合。四角の枠内のラベルが系列を表す。

することを把握することができる。なぜなら、「借りた<sup>6</sup>」にはラベルIが対応し、一方「買った<sup>10</sup>」にはラベルSが対応している。そして、ラベルIはラベルSよりCPに相対的に近いと定義されているためである。ラベル系列と入力文の単語を用いることで、系列モデルは「買った<sup>10</sup>」のSPCよりCPに近い位置に「借りた」が存在することを確率の計算に反映させることができる。

我々は、CRF<sup>[14]</sup>を基にした系列ラベリング技術を用いて、SPに対応するラベル系列を識別する。我々のタスクとCRFのタスクとは2つの違いがある。違いの1つは、CRFは全てのラベル候補からラベル系列を識別するが、我々のラベル系列は、CPはラベルC、SPCはラベルS、それらの間にはラベルIと制約を与えている。もう1つの違いは、CRFは同じ系列データに対するラベル系列を識別することを目的としているが、我々は長さの異なる系列に対するラベル系列を識別する。すなわち、CPとSPCの外側にはラベルを付与しない。この違いはCRFのタスクにおいて、CPとSPCの外側に別のラベル(例えばラベルE)を付与するが、ラベルEに関する素性が全く定義されていない場合に相当する。本稿では、異なる長さの単語列に対応するラベル系列を識別する手法を *partial CRF* と呼ぶ。

我々は *partial CRF* を用いて系列モデルを構築する。系列モデルは、ラベルと系列を利用するようにペアモデルを拡張することで得られる。ラベルを利用するため

に、図3のような CP から SPC の範囲を特定するラベル系列を各 SPC に対して想定する。そして素性にラベルの情報に加えるために、表1の全ての素性テンプレートに素性テンプレート  $\langle l_i, l_j \rangle$  を追加する。ここで  $l_i$  は  $i$  の位置に対応するラベルを表す。例えば、表1の素性テンプレート  $\langle o, s_{i+1}, s_j \rangle$  に  $\langle l_i, l_j \rangle$  を追加することで、素性テンプレートは  $\langle o, s_{i+1}, s_j, l_i, l_j \rangle$  となる。さらに系列を利用するように拡張する。ペアモデルでは、素性テンプレートを適用する際に位置ペア  $(i, j)$  を用いた。それに対して系列モデルでは、素性テンプレートを適用する際に位置ペア  $(i, k)$ ,  $k \in \{k | i < k \leq j \vee j \leq k < i\}$ , および  $(k, j)$ ,  $k \in \{k | i \leq k < j \vee j < k \leq i\}$  を用いる。この際、表1の素性テンプレートでは、2つの位置を表す際に  $i$  と  $j$  を用いているが、素性テンプレートを系列モデルに適用する際には、これらの2つの位置を表す変数のどちらかに  $k$  を用いる。例えば、位置ペア  $(i, k)$  の場合、 $i$  と  $j$  を用いた素性テンプレート  $\langle o, s_{i+1}, s_j, l_i, l_j \rangle$  を適用する際は、 $\langle o, s_{i+1}, s_k, l_i, l_k \rangle$  とする。

系列モデルは SPC の  $j$  が SP である確率を次式で計算する。

$$P(X=j|i, S) = \frac{1}{Z_i} \exp\left(\sum_{k \in MU_{ij}} \mathbf{w}^T \mathbf{f}(i, k, S, o, d, l_i, l_k)\right) + \sum_{k \in MU_{ij}} \mathbf{w}^T \mathbf{f}(k, j, S, o, d, l_k, l_j) \quad (3)$$

ここで、

$$M = \begin{cases} \{m | i < m < j\} & (i < j) \\ \{m | j < m < i\} & (i > j) \end{cases}$$

$$Z_i = \sum_{j \in \{1 \leq j \leq n+1 \wedge j \neq i\}} \exp\left(\sum_{k \in MU_{ij}} \mathbf{w}^T \mathbf{f}(i, k, S, o, d, l_i, l_k)\right) + \sum_{k \in MU_{ij}} \mathbf{w}^T \mathbf{f}(k, j, S, o, d, l_k, l_j)$$

である。 $j$  は SPC の他にラベル系列 ID を表しているため、(3)式はラベル系列を識別していると見なすことができる。

(3)式の  $\exp(\cdot)$  内の最初の項は位置  $i$  と系列内の他の位置との組およびそれらの文脈を利用し、2番目の項は位置  $j$  と系列内の他の位置との組およびそれらの文脈を利用する。

長さの異なる系列を識別するようにモデルを設計す

ることで、系列モデルは距離の影響を自然に扱うことができる。なぜなら、長い系列には多くのラベルがあるため含まれる素性が多く、短い系列にはラベルが少ないため含まれる素性が少ない。この素性の量のバイアスが距離の影響を学習する際に重要な手がかりとなるためである。

### 3.3 訓練データ

提案手法の語順推定モデルの学習には教師有りの訓練データを用いる。このデータはパラレルコーパスと言語間の単語アラインメントから構築する。図4に訓練データの例を示す。対訳文の目的言語側でアラインメントされている単語を左から右にたどっていくと、それらの単語に対応する原言語単語の順番（原言語側の矢印）が教師データとなる。原言語文と原言語側の矢印で表される情報を訓練データとして用いる。

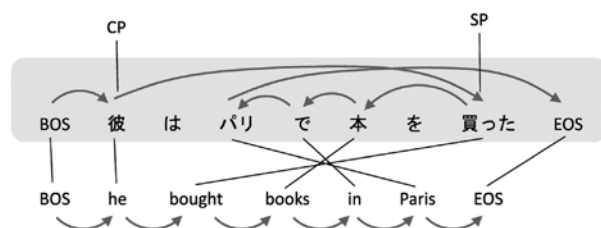


図4 日英翻訳での教師有り訓練データの例。日英の単語間の線は単語アラインメントを表す。

## 4 実験

### 4.1 実験設定

NTCIR-9 特許機械翻訳タスク<sup>[15]</sup>のデータを用いて日英 (JE)、中英 (CE) の特許翻訳の実験を行った。提案手法で用いたフレーズベース SMT の翻訳システムには Moses<sup>[8]</sup> と同等のシステムを用いた。単語アラインメントには GIZA++ を用いた。言語モデルは、対訳データの目的言語側のデータで学習し、5-gram を用いた。翻訳モデルは 40 単語以下の対訳データで学習した。紙面の都合上、詳細な実験設定<sup>[11]</sup>の説明は省略する。

下記の 5 つの語順推定モデルを SMT システムの素性として用いた場合を比較して、提案手法の有効性を評価した。

- ・線形歪みコストモデル<sup>[8]</sup> (LINEAR)

- ・ LINEARとMSD 語彙化語順推定モデル<sup>[9]</sup>(LINEAR+LEX)
- ・ 9クラス分類に基づく語順推定モデル<sup>[12]</sup>(9-CLASS)
- ・ 提案手法のペアモデル (PAIR)
- ・ 提案手法の系列モデル (SEQUENCE)

提案手法の PAIR と SEQUENCE の構築には、計算量を削減するために、翻訳モデルを構築する際に用いたデータのうち 20 万文のデータを用いて学習した。モデルパラメータの推定には L-BFGS<sup>[16]</sup> を用いた。

9-CLASS は、提案手法のモデルで用いたデータと同じ 20 万文を用いて構築した。モデルパラメータの推定には L-BFGS を用いた。

MSD 語彙化語順推定モデルは、翻訳モデルを構築する際に用いたデータを全て用いて構築した。

線形歪みコストモデルは、距離<sup>3</sup>のみに応じてコストが決まるモデルである。

フレーズベース SMT の distortion limit には 10, 20, 30,  $\infty$  の 4 種類の値を用いた。本稿では、開発データを用いてシステム毎に distortion limit を選択した。

また、階層フレーズベース SMT (HIER) の翻訳も行った。このシステムには Moses を用いた。max-chart-span は  $\infty$  に設定した。

## 4.2 評価結果

評価は、BLEU-4 (case insensitive) で行った。図 5 に JE の結果を、図 6 に CE の結果を示す。提案手法の SEQUENCE の結果は、比較したフレーズベース SMT の語順推定モデルの結果よりも評価が高かった。これにより提案手法の系列モデルの有効性が確認された。SEQUENCE の結果は PAIR の結果より評価が高かった。これにより、SEQUENCE で系列を扱うことで候補間の相対的な語順を考慮できることなどが有効であることが確認された。また、階層フレーズベース SMT (HIER) との比較でも SEQUENCE は JE で値が高く、CE でも同等以上の値であった。

提案手法の系列モデルが距離の影響をどれだけ適切に学習できているか調べるために、距離と平均確

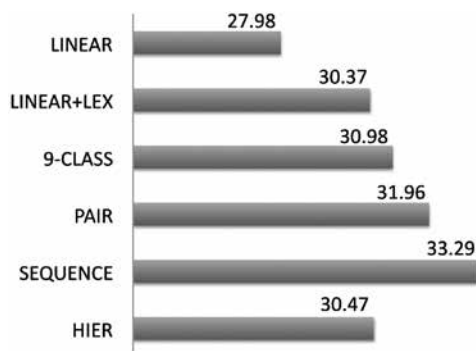


図 5 日英翻訳の評価結果



図 6 中英翻訳の評価結果

率との関係を調べた。具体的には、日英翻訳において、 $\text{distortion} = j - i - 1$  と次の 3 種類の確率との関係を調べた。(1) 日本語テスト文における SEQUENCE の平均確率。(2) 日本語テスト文における PAIR の平均確率。(3) 訓練データにおける最尤推定の確率 (CORPUS)。結果を図 7 に示す。図 7 より、距離クラスが同じ distortion (5 から 20 は同じ距離クラス素性) では、PAIR はほぼ同じ平均確率であるのに対し、SEQUENCE は distortion が大きくなると平均確率が低くなっている。この傾向は CORPUS と同じである。これにより、SEQUENCE が距離による影響を PAIR に比べて適切に学習できていることが確認された。

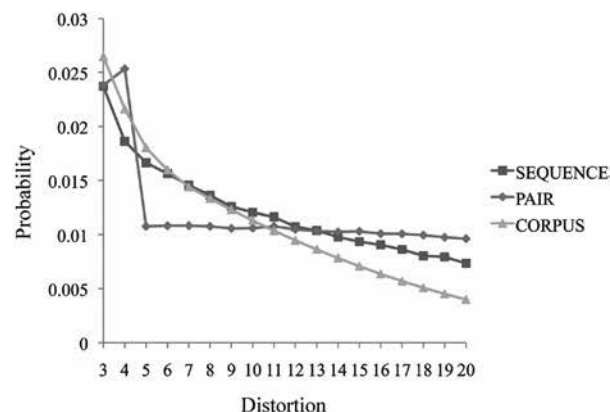


図 7 日英翻訳での distortion と確率との関係

3 より正確には、次節で説明する distortion である。

## 5 まとめ

フレーズベース SMT のための新しい語順推定モデルを提案した。提案した系列モデルは 1 つの確率モデルのみで構成され、候補間の相対的な語順や距離の影響を考慮することができる。日英、中英の特許翻訳の実験で有効性を確認した。

### 参考文献

- [1] Isao Goto, Masao Utiyama, Eiichiro Sumita, Akihiro Tamura, and Sadao Kurohashi. Distortion Model Based on Word Sequence Labeling for Statistical Machine Translation. *ACM Transactions on Asian Language Information Processing* 13, 1, Article 2, 2014.
- [2] Isao Goto, Masao Utiyama, Eiichiro Sumita, Akihiro Tamura, and Sadao Kurohashi. Distortion Model Considering Rich Context for Statistical Machine Translation. In *Proceedings of ACL*, pp. 155-165, 2013.
- [3] Christoph Tillman. A Unigram Orientation Model for Statistical Machine Translation. In *Proceedings of HLT-NAACL*, pp. 101-104, 2004.
- [4] Richard Zens et al. Reordering Constraints for Phrase-Based Statistical Machine Translation. In *Proceedings of Coling*, pp. 205-211, 2004.
- [5] Fei Xia and Michael McCord. Improving a Statistical MT System with Automatically Learned Rewrite Patterns. In *Proceedings of Coling*, pp. 508-514, 2004.
- [6] David Chiang. Hierarchical Phrase-Based Translation. *Computational Linguistics* 33, 2, pp. 201-228, 2007.
- [7] Kenji Yamada and Kevin Knight. A Syntax-based Statistical Translation Model. In *Proceedings of ACL*, pp. 523-530, 2001.
- [8] Philipp Koehn et al. Moses: Open Source Toolkit for Statistical Machine Translation. In *Proceedings of ACL*, pp. 177-180, 2007.
- [9] Philipp Koehn et al. Edinburgh System Description for the 2005 IWSLT Speech Translation Evaluation. In *Proceedings of IWSLT*, 2005.
- [10] Michel Galley and Christopher D. Manning. A Simple and Effective Hierarchical Phrase Reordering Model. In *Proceedings of EMNLP*, pp. 848-856, 2008.
- [11] Richard Zens and Hermann Ney. Discriminative Reordering Models for Statistical Machine Translation. In *Proceedings of WMT*, pp. 55-63, 2006.
- [12] Spence Green et al. Improved Models of Distortion Cost for Statistical Machine Translation. In *Proceedings of HLT-NAACL*, pp. 867-875, 2010.
- [13] Minwei Feng et al. Advancements in Reordering Models for Statistical Machine Translation. In *Proceedings of ACL*, pp. 322-332, 2013.
- [14] John D. Lafferty et al. Conditional Random Fields : Probabilistic Models for Segmenting and Labeling Sequence Data. In *Proceedings of ICML*, pp. 282-289, 2001.
- [15] Isao Goto et al. Overview of the Patent Machine Translation Task at the NTCIR-9 Workshop. In *Proceedings of NTCIR-9*, pp. 559-578, 2011.
- [16] D.C. Liu and J. Nocedal. On the Limited Memory Method for Large Scale Optimization. *Mathematical Programming B* 45, 3, pp. 503-528, 1989.