

# 文頭・文末予測の組み合わせによる文特定

宇田川 拓真 金山 博 吉田 一星

日本アイ・ビー・エム株式会社 東京基礎研究所

Takuma.Udagawa@ibm.com {hkana, issei}@jp.ibm.com

## 概要

文分割とはテキストを連続する文の単位に分割する前処理を指す。本稿ではこれを拡張し、テキスト中の文でないノイズを除去しつつ、解析対象とすべき文のみを抽出する文特定タスクの定式化を行う。また、文頭・文末予測を組み合わせることで文特定を行う新たな手法を提案する。最後に、Universal Dependencies (UD) のアノテーションに基づく言語横断的なベンチマーク手法を設計し、UD の英語コーパス English Web Treebank (EWT) を用いた実験により提案手法の精度・有用性を検証する。

## 1 はじめに

自然言語処理において文は重要な処理単位の一つである [1, 2, 5]。一般的に、テキストはまず文末予測に基づく文分割という前処理によって連続する文の単位 (文単位と呼ぶ) に分割される [4, 3, 11]。しかし、実際のテキストデータは必ずしも綺麗な文単位のみを含むとは限らない。例えば、ウェブデータはメタ情報・文断片・非言語的記号など文と見なせないノイズ (非文単位と呼ぶ) を含み得る。従来の文分割では非文単位もいずれかの文単位に属すると仮定されるが、これらのノイズは構文解析など下流の処理やその評価において障害となることがある [12]。

そこで本稿では文分割の代替として、テキストを連続する文単位と非文単位に分割する文特定タスクの定式化を行う。それと同時に、文頭と文末予測を組み合わせることで文特定を行う新たな手法を提案する。表 1 に示すように、提案手法では文頭・文末予測によって文単位を特定しつつ、それ以外を非文単位と見なして区別・除去することができる。

次に、文特定タスクにそのまま適用可能なベンチマークが存在しないため、Universal Dependencies (UD) [8] のアノテーションを基に文特定のベンチマークを設計する。具体的には、UD が提供する文境界と述語項構造を活用して、任意の UD コーパスか

ら文特定の正解データを作成する手法を紹介する。文単位・非文単位の分類基準を目的に応じて柔軟に調整可能であることもこの手法の特長である。

最後に、UD の英語コーパスとして広く使われている English Web Treebank (EWT) [10] を用いて実験を行う。これにより、提案手法は高い精度で文単位を特定できること、また正確な文単位の特特定が構文解析等の下流タスクにおいて有用であることを示す。

## 2 タスク設計・手法

本節ではまず、文分割のタスク設計と手法を再定義する。次に、文分割の考えを拡張することで文特定のタスク設計と手法を導出できることを示す。

### 2.1 文分割

入力テキストを  $\mathbf{W} = (w_0, w_1, \dots, w_{N-1})$  とし、各  $w_i$  は単語とする (文字・サブワードでも可)。また、テキストのスパン  $\mathbf{W}[i:j] = (w_i, \dots, w_{j-1})$  およびその連結  $\mathbf{W}[i:j] \oplus \mathbf{W}[j:k] = \mathbf{W}[i:k]$  を定義する。

文分割では入力テキストを連続する  $M$  個の文単位に分割する。<sup>1)</sup> つまり、文単位の境界のインデックスを  $\mathbf{B} = (b_0, b_1, \dots, b_M)$  (ただし  $b_0 = 0, b_M = N$ ) とすると、 $\bigoplus_{i=1}^M \mathbf{W}[b_{i-1}:b_i] = \mathbf{W}$  となる。ここで、スパン  $\mathbf{W}[i:j]$  が一つの文単位となる確率を  $p_{\text{SU}}(\mathbf{W}[i:j])$  とおくと、最適な文分割は以下の解  $\mathbf{B}$  を求めるタスクとして定式化できる：

$$\arg \max_{\mathbf{B}} \prod_{i=1}^M p_{\text{SU}}(\mathbf{W}[b_{i-1}:b_i]) \quad (1)$$

次に、 $w_i$  が文末となる確率のモデル  $p_{\text{EOS}}(w_i|\mathbf{W};\theta)$  を導入する。<sup>2)</sup> 通常、モデルのパラメータ  $\theta$  は文末 (文境界) の正解データを用いて事前学習され [11]、このモデルを用いて文単位の確率  $p_{\text{SU}}(\mathbf{W}[i:j])$  を以下のように定義できる：

$$p_{\text{SU}}(\mathbf{W}[i:j]) = p_{\text{EOS}}(w_{j-1}) \prod_{i \leq k < j-1} (1 - p_{\text{EOS}}(w_k)) \quad (2)$$

1)  $M$  は変数であることに注意。

2) 以降  $\mathbf{W}, \theta$  を省略して  $p_{\text{EOS}}(w_i)$  のように表記する。

表 1 文分割と文特定の比較. 文分割では文末予測 (E) によってテキストを連続する文単位 (青いスパン) に分割する. 文特定では文頭予測 (B) と文末予測 (E) の間のスパンを文単位として特定し, それ以外のスパンを非文単位と見なす.

入力テキスト	Thank you. - TEXT.htm << File: TEXT.htm >> I was thinking of converting it to a hover vehicle.			
(EWT の例)	I might just sell the car and get you to drive me around all winter.			
文分割	E			E
	Thank you. - TEXT.htm << File: TEXT.htm >> I was thinking of converting it to a hover vehicle.			
		E		
	I might just sell the car and get you to drive me around all winter.			
文特定	B	E	B	E
	Thank you. - TEXT.htm << File: TEXT.htm >>		I was thinking of converting it to a hover vehicle.	
	B			E
	I might just sell the car and get you to drive me around all winter.			

つまり  $W[i:j]$  が文単位となる確率は, そこに含まれる単語のうち最後の単語  $w_{j-1}$  のみが文末となる確率に等しい. この式 (2) を式 (1) に代入すると,

$$\begin{aligned}
 (1) &= \arg \max_{\mathbf{B}} \sum_{i=1}^M \left\{ \log p_{\text{EOS}}(w_{b_i-1}) + \sum_{b_{i-1} \leq j < b_i-1} \log(1 - p_{\text{EOS}}(w_j)) \right\} \\
 &= \arg \max_{\mathbf{B}} \sum_{i \in \mathbf{B}_{\text{EOS}}} \log p_{\text{EOS}}(w_i) + \sum_{i \notin \mathbf{B}_{\text{EOS}}} \log(1 - p_{\text{EOS}}(w_i))
 \end{aligned} \quad (3)$$

ここで  $\mathbf{B}_{\text{EOS}} = \{b_i - 1 \mid i \in (1, 2, \dots, M)\}$  は  $\mathbf{B}$  が定義する文末のインデックスを指す. 式 (3) は自明な最適化問題であり,  $\mathbf{B}_{\text{EOS}} = \{i \in (0, 1, \dots, N-1) \mid p_{\text{EOS}}(w_i) \geq 0.5\}$  となる  $\mathbf{B}$  が最適解となる.

## 2.2 文特定

文特定では入力テキストを連続する文単位または非文単位に分割する. そのため,  $a_i$  を  $W[b_{i-1}:b_i]$  が文単位の場合 1, 非文単位の場合 0 を取る変数とすると, 最適な文特定は以下の解  $\mathbf{A} = (a_1, a_2, \dots, a_M)$  および  $\mathbf{B}$  を求めるタスクとして定式化できる:

$$\arg \max_{\mathbf{B}, \mathbf{A}} \prod_{i=1}^M p_{\text{SU}}(\mathbf{W}[b_{i-1}:b_i])^{a_i} p_{\text{NSU}}(\mathbf{W}[b_{i-1}:b_i])^{1-a_i} \quad (4)$$

ここで  $p_{\text{SU}}(\mathbf{W}[i:j])$ ,  $p_{\text{NSU}}(\mathbf{W}[i:j])$  はそれぞれ  $W[i:j]$  が一つの文単位, 非文単位となる確率を指す. 前項の考えを拡張すると, これらの確率は事前学習された文頭確率のモデル  $p_{\text{BOS}}(w_i)$  と文末確率のモデル  $p_{\text{EOS}}(w_i)$  を用いて以下のように定義できる:

$$\begin{aligned}
 p_{\text{SU}}(\mathbf{W}[i:j]) &= p_{\text{BOS}}(w_i) \prod_{i < k \leq j-1} (1 - p_{\text{BOS}}(w_k)) \\
 &\quad \times p_{\text{EOS}}(w_{j-1}) \prod_{i \leq k < j-1} (1 - p_{\text{EOS}}(w_k)) \\
 p_{\text{NSU}}(\mathbf{W}[i:j]) &= \prod_{i \leq k \leq j-1} (1 - p_{\text{BOS}}(w_k)) \times \prod_{i \leq k \leq j-1} (1 - p_{\text{EOS}}(w_k))
 \end{aligned} \quad (5)$$

つまり  $W[i:j]$  が文単位となる確率は, 最初の単語  $w_i$  のみが文頭となり, 最後の単語  $w_{j-1}$  のみが文末となる確率に等しい. また,  $W[i:j]$  が非文単位となる確率は文頭および文末が一切含まれない確率に等しい.<sup>3)</sup> これらの式 (5) を式 (4) に代入すると,

$$\begin{aligned}
 (4) &= \arg \max_{\mathbf{B}, \mathbf{A}} \sum_{i=1}^M \left\{ a_i \log p_{\text{SU}}(\mathbf{W}[b_{i-1}:b_i]) \right. \\
 &\quad \left. + (1 - a_i) \log p_{\text{NSU}}(\mathbf{W}[b_{i-1}:b_i]) \right\} \\
 &= \arg \max_{\mathbf{B}, \mathbf{A}} \sum_{i \in \mathbf{B}_{\text{AOS}}} \log p_{\text{BOS}}(w_i) + \sum_{i \notin \mathbf{B}_{\text{AOS}}} \log(1 - p_{\text{BOS}}(w_i)) \\
 &\quad + \sum_{i \in \mathbf{B}_{\text{EOS}}} \log p_{\text{EOS}}(w_i) + \sum_{i \notin \mathbf{B}_{\text{EOS}}} \log(1 - p_{\text{EOS}}(w_i))
 \end{aligned} \quad (6)$$

ここで  $\mathbf{B}_{\text{AOS}}^{\mathbf{A}} = \{b_{i-1} \mid i \in (1, 2, \dots, M), a_i = 1\}$  は文頭のインデックス,  $\mathbf{B}_{\text{EOS}}^{\mathbf{A}} = \{b_i - 1 \mid i \in (1, 2, \dots, M), a_i = 1\}$  は文末のインデックスを指す. 付録 A に示すように, 式 (6) の最適解  $\mathbf{B}, \mathbf{A}$  は動的計画法を用いて厳密かつ効率的に求めることができる.

## 3 評価

本節では, 任意の言語の UD コーパスから文特定のベンチマークを構築する手法について説明する. 具体的には UD のアノテーションを基に, まず (1) 文単位・非文単位の境界を決定し, 次に (2) 各区分を文単位または非文単位に分類する.

まず手順 (1) では, UD における文境界<sup>4)</sup> を文単位・非文単位の境界として採用する. 手順 (2) においては, 言語学における Lexical Sentence [9] の考え方を採用し, 単語間の依存関係を基に文単位を定義する. 具

3) この定義に基づくと, 任意の  $W[i:j] \oplus W[j:k] = W[i:k]$  に対して  $p_{\text{NSU}}(\mathbf{W}[i:k]) = p_{\text{NSU}}(\mathbf{W}[i:j]) \times p_{\text{NSU}}(\mathbf{W}[j:k])$  が成立する. つまり, 連続する非文単位の区別は行われない.

4) UD において, 文境界における「文」とは表 2 のような非文単位を含む広義の意味で使われている.



表 4 文単位の抽出タスクの実験結果. 単語レベルの BIO ラベルとスパンレベルの F1 スコアによって評価する.

手法	EWT テスト ( $p_{cc} = 0.5$ )				EWT テスト ( $p_{cc} = 0$ )				EWT テスト (後処理)			
	B	I	O	Span	B	I	O	Span	B	I	O	Span
文末による文分割	85.6	97.3	66.6	72.8	78.0	95.1	6.0	58.2	90.2	97.5	71.3	81.6
文末による文分割 (+文末強制)	79.8	95.9	0.0	60.4	77.8	95.1	0.0	57.7	81.7	95.7	0.0	62.3
文頭・文末による文特定	<b>94.3</b>	<b>98.7</b>	<b>86.1</b>	<b>87.3</b>	<b>93.0</b>	<b>98.2</b>	<b>81.7</b>	<b>84.1</b>	<b>94.7</b>	<b>98.4</b>	<b>83.9</b>	<b>88.8</b>

表 5 下流タスク (単語区切り/ Words・品詞タグ付け/ UPOS・依存構造解析/ LAS) の評価結果.

EWT 訓練・検証データ	EWT テストデータ								
	全体			文単位のみ			非文単位のみ		
	Words	UPOS	LAS	Words	UPOS	LAS	Words	UPOS	LAS
全体	98.6	96.2	89.7	98.8	96.5	92.0	96.5	92.0	81.4
文単位のみ	-	95.7	89.0	-	96.5	90.2	-	88.3	75.4
非文単位のみ	-	72.6	36.1	-	70.7	31.8	-	91.0	77.1

見られる<sup>9)</sup>一方で, 文特定の手法ではこれらに対しても高い性能を維持していることが分かる.

以上の結果より, 非文単位を除いて文単位を特定するには文末の予測のみ (= 文分割) では不十分であり, 文頭・文末の予測を組み合わせること (= 文特定) が重要であると言える.

### 4.3 下流タスクの評価

最後に, 文特定の下流タスクにおける有用性について, 単語区切り・品詞タグ付け・依存構造解析の3つに焦点を当てて検証する.

まず背景として, EWT は非文単位のようなノイズを含むため, 解析器の訓練・評価の両面で支障があることが知られている [12]. そこで本稿では, 理想的な文単位・非文単位の区分<sup>10)</sup>に基づいて解析器を訓練・評価することの利点について調査する.

下流タスクの解析器には Trankit [7] を用いる. 手順としては UD-EWT v2.6 の全体 (文単位と非文単位) を用いて単語区切りを訓練し<sup>11)</sup>, 続いて UD-EWT v2.10 の全体・文単位のみ・非文単位のみそれぞれで品詞タグ付けと依存構造解析を訓練する. 最後に, UD-EWT v2.10 の全体・文単位のみ・非文単位のみそれぞれで各解析器を評価する.

表 5 に下流タスクの評価結果を示す. まず EWT 全体で解析器を訓練した場合, 全ての下流タスクにおいて文単位のみで評価したスコアが高く, 非文単位

のみで評価したスコアが顕著に低くなっている. このことから, 文単位は予測に適した綺麗なデータが多く, 逆に非文単位は予測しづらいイレギュラーなデータが多いことが分かる. これらを同一視せずに区別して評価を行うことは, 解析器の特徴を正しく理解するために重要と考えられる.

また, 文単位のみで解析器を訓練した場合, 学習データは少なくなるにもかかわらず, 全体で訓練した場合と近い性能が出ることが分かる.<sup>12)</sup> 例えば文単位に対する品詞タグ付けでは, 全体で訓練した場合と同等の性能が出ている. 逆に非文単位のみで訓練した場合の性能は全体的に低く, 訓練データとしても有用でないノイズが多いと考えられる.<sup>13)</sup> 今後の応用として, 例えば文単位・非文単位の分類基準 (3 節) を調整することで訓練に最適なデータのフィルタリングを行えると考えられる.

以上の結果より, 正確な文特定は解析器の訓練・評価の両面で有用であると考えられる.

## 5 まとめ

本稿では文分割の考え方を拡張し, 非文単位を除きつつ文単位を特定する文特定のタスクと手法を提案した. また, EWT を用いた実験により文特定の精度と有用性を確認した. 今後の課題として, 文特定のさらなる精度向上, 異なる言語・ドメインにおける有効性の検証, 実際の (誤りを含む) 予測に基づく下流タスクの評価・分析などが挙げられる.

9) 文分割では予測できる非文単位 (O ラベル) の割合が入力が長くなるほど小さくなり, 結果的に性能が大幅に劣化する.

10) つまり 3 節で示した基準で分類した正解の区分を用いる. 現実的には文特定の予測の段階で誤りが生じる可能性があるが, その点まで考慮した分析は今後の課題とする.

11) 複数語トークンによるエラーを回避するため, それを持たない旧バージョンで訓練した単語区切りモデルで固定した.

12) ただし非文単位のみで評価した結果は異なる傾向を持ち, やはりイレギュラーなデータを多く含むことを示唆している.

13) フェアな比較のためには訓練に使うデータ量を揃える必要があるが, 付録 C に示すようにこの設定でも同様の傾向が確認された.

## 参考文献

- [1] Timothy Dozat and Christopher D Manning. Deep biaffine attention for neural dependency parsing. In *Proc. of ICLR*, 2017.
- [2] Timothy Dozat and Christopher D. Manning. Simpler but more accurate semantic dependency parsing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 484–490, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [3] Dan Gillick. Sentence boundary detection and the problem with the U.S. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers*, pp. 241–244, Boulder, Colorado, June 2009. Association for Computational Linguistics.
- [4] Tibor Kiss and Jan Strunk. Unsupervised multilingual sentence boundary detection. *Computational Linguistics*, Vol. 32, No. 4, pp. 485–525, 2006.
- [5] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, Vol. 8, pp. 726–742, 2020.
- [6] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [7] Minh Van Nguyen, Viet Dac Lai, Amir Pouran Ben Veyseh, and Thien Huu Nguyen. Trankit: A light-weight transformer-based toolkit for multilingual natural language processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pp. 80–90, Online, April 2021. Association for Computational Linguistics.
- [8] Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. Universal Dependencies v2: An evergrowing multilingual treebank collection. In *Proceedings of the 12th Language Resources and Evaluation Conference (LREC 2020)*, pp. 4034–4043, Marseille, France, May 2020. European Language Resources Association.
- [9] Geoffrey Nunberg. *The linguistics of punctuation*. No. 18. Center for the Study of Language (CSLI), 1990.
- [10] Natalia Silveira, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Chris Manning. A gold standard dependency corpus for English. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pp. 2897–2904, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA).
- [11] Rachel Wicks and Matt Post. A unified approach to sentence segmentation of punctuated text in many languages. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 3995–4007, Online, August 2021. Association for Computational Linguistics.
- [12] 金山博, 大湖卓也. UD\_English-EWT の付き合い方. 言語処理学会第 28 回年次大会予稿集, March 2022.

表 6 English Web Treebank (EWT) に含まれる文単位の代表例. 各行は一つの文単位に対応する.

President Bush on Tuesday nominated two individuals to replace retiring jurists on federal courts in the Washington area. Unfortunately, Mr. Lay will be in San Jose, CA participating in a conference, where he is a speaker, on June 14.  
 “In 1972, there was an enormous glut of pilots,” Campenni says.  
 PS – There is a happy hour tonight at Scudeiros on Dallas Street (just west of the Met Garage) beginning around 5:00.  
 2) Your vet would not prescribe them if they didn't think it would be helpful.  
 BUT EVERYONE HAS THERE OWN WAY!!!!!!  
 The motel is very well maintained, and the managers are so accomodating, it's kind of like visiting family each year! :-)  
 where can I find the best tours to the Mekong Delta at reasonable prices?

表 7 訓練データ量を揃えた場合の下流タスク (単語区切り・品詞タグ付け・依存構造解析) の評価結果.

EWT 訓練・検証データ	EWT テストデータ								
	全体			文単位のみ			非文単位のみ		
	Words	UPOS	LAS	Words	UPOS	LAS	Words	UPOS	LAS
全体	98.6	94.6	82.8	98.8	95.1	83.7	96.5	90.2	71.4
文単位のみ	-	94.0	82.4	-	94.9	83.9	-	85.5	65.3
非文単位のみ	-	72.6	36.1	-	70.7	31.8	-	91.0	77.1

## A 動的計画法

2.2 項の式 (6) の最適解は動的計画法を用いて求めることができる. 具体的には, 任意の  $k \leq N-1$  に対して  $\mathbf{W}^{\leq k} = (w_0, \dots, w_k)$  までの最適な文頭・文末予測を考える. つまり, 式 (6) における  $\mathbf{W}$  を  $\mathbf{W}^{\leq k}$  に置き換えた最適化問題・目的関数を考える.

$\mathbf{W}^{\leq k}$  までの予測は不完全であるため, 文単位の途中で終わる (= 最後の予測が文頭である) 場合と, 文単位の外で終わる (= 最後の予測が文頭でない) 場合の 2 通りが考えられる.  $\log p_{\text{IS}}(k+1)$  を前者の場合の目的関数の最大値,  $\log p_{\text{OS}}(k+1)$  を後者の場合の目的関数の最大値と置く. すると, 予測は必ず文単位の外から始まるので  $\log p_{\text{IS}}(0) = \log 0 = -\infty$ ,  $\log p_{\text{OS}}(0) = \log 1 = 0$  と初期化でき, これらの値を以下のように更新して求めることができる.

$$\begin{aligned} \log p'_{\text{IS}}(i) &= \max \{ \log p_{\text{IS}}(i) + \log (1 - p_{\text{BOS}}(w_i)), \\ &\quad \log p_{\text{OS}}(i) + \log p_{\text{BOS}}(w_i) \} \\ \log p'_{\text{OS}}(i) &= \log p_{\text{OS}}(i) + \log (1 - p_{\text{BOS}}(w_i)) \\ \log p_{\text{IS}}(i+1) &= \log p'_{\text{IS}}(i) + \log (1 - p_{\text{EOS}}(w_i)) \\ \log p_{\text{OS}}(i+1) &= \max \{ \log p'_{\text{IS}}(i) + \log p_{\text{EOS}}(w_i), \\ &\quad \log p'_{\text{OS}}(i) + \log (1 - p_{\text{EOS}}(w_i)) \} \end{aligned} \quad (7)$$

この更新式ではまず, 文頭確率  $p_{\text{BOS}}(w_i)$  に基づいて  $p_{\text{IS}}(i) \rightarrow p'_{\text{IS}}(i)$ ,  $p_{\text{OS}}(i) \rightarrow p'_{\text{OS}}(i)$  の更新を行う. 次に, 文末確率  $p_{\text{EOS}}(w_i)$  に基づいて  $p'_{\text{IS}}(i) \rightarrow p_{\text{IS}}(i+1)$ ,  $p'_{\text{OS}}(i) \rightarrow p_{\text{OS}}(i+1)$  の更新を行う.

最終的に予測は文単位の外で終わる必要があるた

め,  $\log p_{\text{OS}}(N)$  が式 (6) の目的関数の最大値となる. 式 (6) の最適解  $B, A$  は, 更新式 (7) をバックトラックすることで簡単に求めることができる.

## B 文単位の例

表 6 に EWT に含まれる文単位の代表例を示す. これらの例に示されるように, 本稿の基準では節を成す述語項構造を持つ区分のみが解析対象とすべき文単位として判定される.

## C 下流タスク評価の追加結果

表 7 に訓練データ量を揃えた場合の下流タスクの評価結果を示す. 手順としては, 非文単位のデータ量が最も少ないため, 訓練データに含まれる単語数を非文単位におおよそ合わせて訓練を行った.<sup>14)</sup>

この結果から, やはり文単位のみで解析器を訓練した場合, 全体で訓練した場合と同等の性能が出ることが分かる. また, データ量を揃えた場合でも非文単位のみで訓練した場合の性能は全体として低いことが分かる. このことから, 下流タスクの学習には文単位が最も有用であり, 非文単位の貢献は非常に小さいことが分かる. ただし非文単位のみで評価した場合の結果は異なる傾向を示しており, 非文単位に含まれるデータは (アノテーションが一貫していないなど) イレギュラーであることが示唆される.

14) 訓練データの単語数をおおよそ合わせるために, 全体で訓練する場合はランダムに 622 例, 文単位のみで訓練する場合はランダムに 518 例, 非文単位のみで訓練する場合は 2,187 例 (全ての例) を使用した.