

# 盛り上がり時間帯におけるツイートの言語的特性の解析

藤沼 祥成<sup>1,2,a)</sup> 横野 光<sup>2,b)</sup> Pascual Martínez-Gómez<sup>2,3,c)</sup> 相澤 彰子<sup>1,2,d)</sup>

**概要:** あるイベントの盛り上がりに対して、それに関するツイートにも変化が現れその変化に着目することで盛り上がりを検出することが可能であると考えられる。本研究ではこの盛り上がり時間帯中のツイートに用いられている表現の特性を解析することを試みる。はじめに各時間帯のツイート集合とツイートより構築した言語モデルの関係をクロスエントロピーで算出した。実験結果より複数のハッシュタグ間における一部の盛り上がり時間帯のツイートはツイートより構築した n-gram 言語モデルに従うことを示す。また、盛り上がっている時間帯とそうでない時間帯において、クロスエントロピーにおいて統計的に有意差があることを示した ( $p < 0.02$ )。また、n-gram 言語モデルでは捉えられない素性も検討するため、Support Vector Machine (SVM) と Random Forest により各ツイートを盛り上がり時間帯の二値分類を行い、盛り上がり時間帯の特徴として漢字数が少ないことが明らかになった。

## 1. はじめに

近年マイクロブログの発展により、ソーシャルメディア上でユーザが、実際に遭遇したイベントなどに関してリアルタイムに言及することが増えてきた。2013年現在では日に5億ものツイートが投稿され<sup>\*1</sup>、これらのツイートの中にはイベント中の盛り上がりによって発言されるツイートも含まれる。このようなイベントの盛り上がりはTwitterではツイート数の急激な増加によって現れることも多く、2010年のワールドカップにおける日本対カメルーン戦においては、日本のゴールの瞬間に一秒あたり2,940ツイートという当時の一秒あたりの最高ツイート数を記録した<sup>\*2</sup>。図1に日本対カメルーン戦における一分あたりのツイート数の変化とツイート数が最も多くみられた上位3つの盛り上がり時間帯に何が起きたかを示す。

このようなイベントにおけるTwitter上の盛り上がりを検出できれば、実世界での重要な瞬間を捜す手がかりになる。それを利用した研究として、例えばLanaganら[11]は、

スポーツの試合に関する要約の自動生成において、Twitter上の盛り上がり時間帯を利用して重要な瞬間を捉える手法を提案している。

スポーツに限らず、Twitter上のイベントにおける盛り上がりを検出する上で時間当たりのツイート数に着目した盛り上がりの指標が提案されている[8], [10], [19]。また、バースト検出[7]、Topic Detection and Tracking[12], [15]を含むTwitter上のタスクにおいては単位時間当たりのツイート数に焦点が当てられている。

ここで、従来の研究では盛り上がっている時間帯の言語的特性は注目されていない。しかし、表1にあるように、盛り上がっているときのツイートと平常時でのツイートには、表現的な違いが見られる。これは盛り上がっているときにはユーザは興奮していることが多く、その感情がツイートにも反映されているためと考えられる。ツイート特有の言語的特性に関しては、Brodyら[2]がTwitter上で文字の繰り返しが起こる単語は極性辞書に含まれる単語である割合が高いことを示しているが、実世界とのイベントとは関連付けられていない。

そこで、本稿では盛り上がっている時間帯における日本語のツイートの言語的特性を解析する。具体的には日本語ツイート全体よりサンプリングしたツイートをもとに言語モデルを構築し、盛り上がり時間帯のツイートと比較する。以下2章にて関連研究を、3章にて使用した3つのデータセットについて述べる。4章にてツイート言語モデルとクロスエントロピーに関して述べ、5章では各ツイートにおける盛り上がり時間帯のツイートかどうかの二値分類を行い、その特徴量の解析を行う。6章で結論を述べる。

<sup>1</sup> 東京大学  
東京都文京区本郷 7-3-1

<sup>2</sup> 国立情報学研究所  
東京都千代田区一ツ橋 2-1-2

<sup>3</sup> お茶の水女子大学  
東京都文京区大塚 2-1-1

a) y\_fujinuma@nii.ac.jp

b) yokono@nii.ac.jp

c) pascual@nii.ac.jp

d) aizawa@nii.ac.jp

\*1 <https://blog.twitter.com/2013/new-tweets-per-second-record-and-how>

\*2 <https://blog.twitter.com/2010/big-goals-big-game-big-records>



表 2 6つのハッシュタグの統計とそのターゲット区間. 時刻は全て UTC +0 である.

Table 2 Statistics of six hashtags and its target interval. All time are UTC +0.

ハッシュタグ	ジャンル	日時	ツイート数
#aibou	テレビドラマ	2012-03-21T10:31 - 2012-03-21T14:06	20,681
#hanshin	野球	2012-04-20T08:39 - 2012-04-20T12:54	6,176
#ACV	オンラインゲーム	2012-02-13T12:08 - 2012-02-13T15:41	1,562
#agqr	ラジオ番組	2012-02-15T11:39 - 2012-02-15T14:08	13,434
#figureskate	フィギュアスケート	2012-04-20T10:00 - 2012-04-20T12:20	1,410
#momoclo	音楽	2012-02-11T15:36 - 2012-02-11T17:21	1,823

ザからの反応である. そのため用いたツイートデータに関しては前処理としてボットからのツイートとリツイートを意味する文字列 ‘RT’ を含むツイートを削除した. ボットのフィルタリングでは, Twitter クライアントの情報を利用した. ツイート数が上位のクライアントにはボットが含まれにくいため, サンプリングした3日間におけるツイートに対して, クライアントから人の発言で用いられているものを人手で上位43クライアントを選択し, それらからの発言のみを利用した. このクライアントからのツイートはサンプリングした3日間全体のツイートの90%以上を含んでいる.

以上のフィルタリングを施したデータより言語モデル構築用, クロスエントロピー算出用, 特徴量解析用の3つのサブデータセットを構築した.

### 3.1 言語モデル構築用データセット

日本語ツイートの言語モデル (以下, ツイート言語モデルと呼ぶ) の構築に2012年2月12日中の全日本語ツイート, 合計413,008,939ツイートよりサンプリングした50,000ツイートをを用いた.

またツイート言語モデルとは別にハッシュタグ#aibouを含むツイートからなる言語モデル (以下, #aibou言語モデルと呼ぶ) を100ツイートを超える時間帯のみからの計15,663ツイートより構築した.

### 3.2 クロスエントロピー算出用データセット

ハッシュタグのジャンルに関わらない, 盛り上がり時間帯における共通の性質を解析するため, 様々なジャンルのハッシュタグのツイートを収集した. ハッシュタグにおいて盛り上がりは実世界のイベントが起きている最中に起こりやすいため, その時間帯をハッシュタグにおけるターゲット区間とする. ターゲット区間は中央値の10倍を目安に設定し, ツイートを抽出した. また, 盛り上がり時間帯以外の時間帯のツイートも含めるため, ターゲット区間の前後20分の時間におけるツイートもデータセットに含めた. 表2に実験に用いたデータセットの詳細を示す.

### 3.3 特徴量解析用データセット

各ツイートが盛り上がり時間帯のものか, その他の時間帯のものかどうかを判定する二値分類問題を考え, これに有効な特徴量から盛り上がり時の言語的特性を解析する. この実験のために150ツイートずつ盛り上がりしている時間帯とそうでない時間帯より抽出した, 偏りのないデータセットを構築した. 具体的な抽出手順を以下に示す.

- (1) 手動で設定したターゲット区間を3分割する.
- (2) 各分割においてツイート数が上位の時間帯より, ツイートをテストデータとしてツイート数の合計が50件となるまで取得していく.
- (3) 最後の時間帯中のツイートを合わせて, サンプリングしたツイートが50件を超える場合は50件になるようランダムサンプリングする.

## 4. 言語モデルによる解析

はじめにツイート言語モデルを構築し, 各時間帯に対応するツイート集合とのクロスエントロピーを算出する.

### 4.1 言語モデルの構築

文字数  $l$  のツイート  $t$ , ツイート中の文字  $t_i$  に対し,  $t$  の確率を以下のように定義する.

$$P(t) = \prod_{i=1}^l P(t_i | t_{i-1}, \dots, t_{i-n+1}).$$

CMU-Cambridge statistical language model toolkit [5] を用いて言語モデルを構築し, Katz バックオフスムージング [9] を用いた.

また, Neubig ら [14] が用いた手法であり, かつウェブ上の日本語において形態素解析の精度が十分でない [17], [21] ため, 文字 7-gram モデルを採用した.

### 4.2 クロスエントロピー

Danescu-Niculescu-Mizil ら [6] の研究に準じて, あるツイート集合  $T$  に対する言語モデルからの距離を数値化するため, クロスエントロピーを採用した.

クロスエントロピー  $H$  を以下に示す.

$$H(T) = -\frac{1}{N} \sum_i^N \log P(T_i)$$

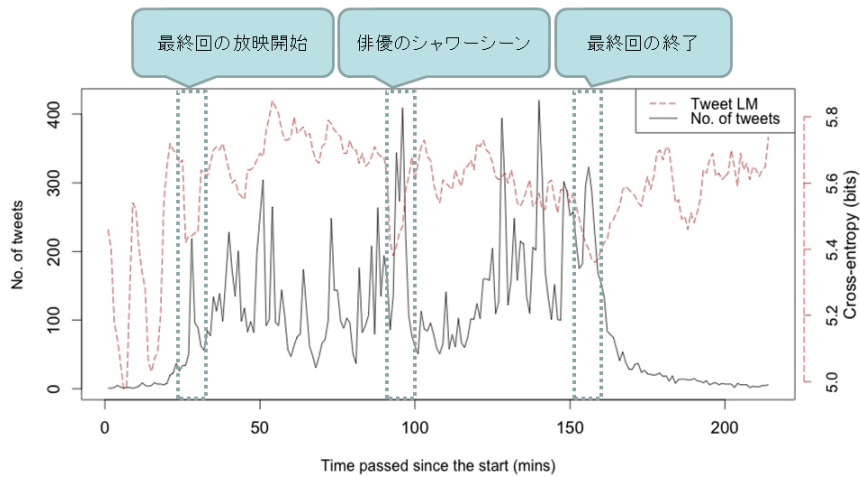


図 2 #aibou におけるツイート数の変化とツイート言語モデルを用いて算出したクロスエントロピーの変化.

Fig. 2 Cross-entropy between the general tweet language model and #aibou data.

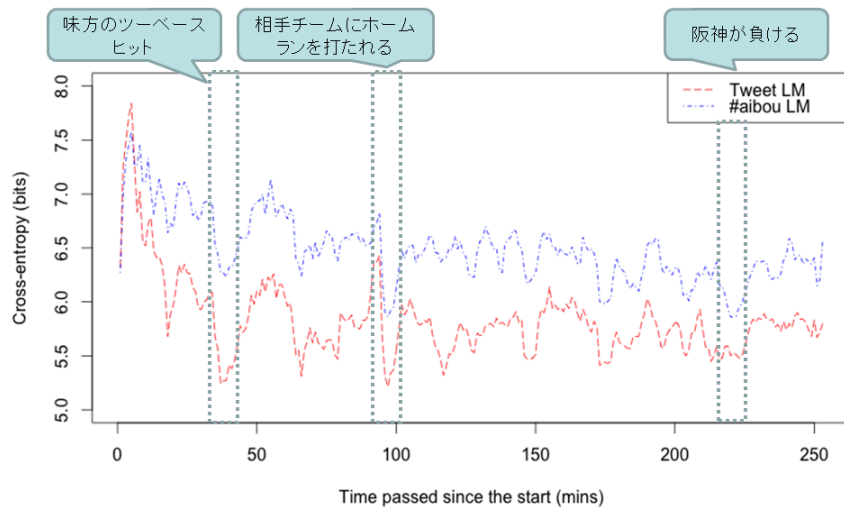


図 3 #hanshin データに対し、ツイート言語モデルと#aibou 言語モデルをそれぞれ用いてクロスエントロピーを算出した結果.

Fig. 3 Cross-entropy of #hanshin data compared with the general tweet language model and the #aibou language model.

なおクロスエントロピー算出時に前後 2 分間のツイートを含めた。これはクロスエントロピーを算出する際サンプリングエラーに弱いことが知られており [3], サンプリングエラーに対して頑健にするためである。

#### 4.3 クロスエントロピーの算出結果

各データに対するクロスエントロピーの算出結果を図 2, 図 3, 図 4, 図 5 に示す。

イベントの始まりと終わりの時間帯はクロスエントロピーが低くなることわかる。図 2 では#aibou において「はじめた」等、イベントの始まりを示す多くツイートが多く見られた。

より詳細に解析するため、ハッシュタグ#aibou の盛り上がり時間帯のツイートだけからなる#aibou 言語モデル

を構築した。図 3 は#hanshin に#aibou 言語モデルとツイート言語モデルをそれぞれ適用した結果を示す。#aibou 言語モデルとツイート言語モデル共にクロスエントロピーの変動が似たような傾向を示している。#aibou 言語モデルを適用した際にクロスエントロピーが最も低下した瞬間は 1) 相手チームがホームランを打った瞬間, 2) 味方チームが負けた瞬間, である。一方ツイート言語モデルでは 1) 相手チームがホームランを打った瞬間, 2) 味方チームがツーベースヒットを打った瞬間, である。#hanshin においては#aibou 言語モデルはより重要な場面においてクロスエントロピーが低下している。

しかし#aibou 言語モデルより算出したクロスエントロピーではツイート言語モデルで算出したクロスエントロピーと比較した結果、ツイート言語モデルのクロスエント

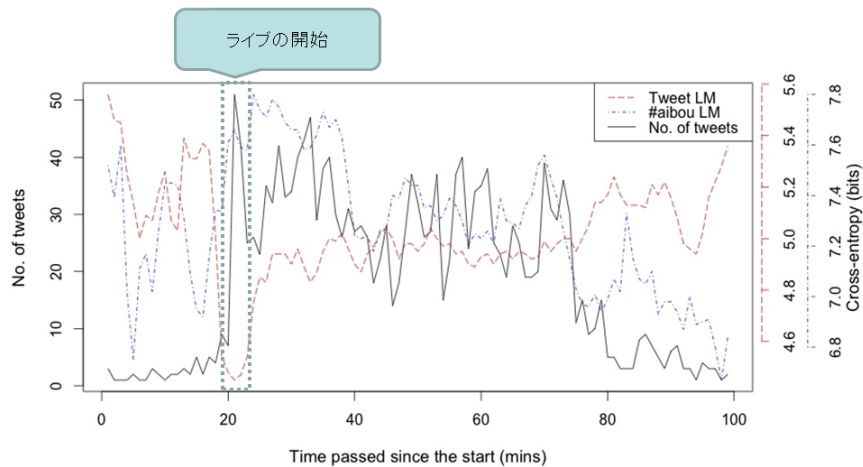


図 4 #momoclo データに対し、ツイート言語モデルと#aibou 言語モデルをそれぞれ用いてクロスエントロピーを算出した結果.

Fig. 4 The occurrence of #momoclo data shown together with the cross-entropy between the general tweet language model and the #aibou language model.

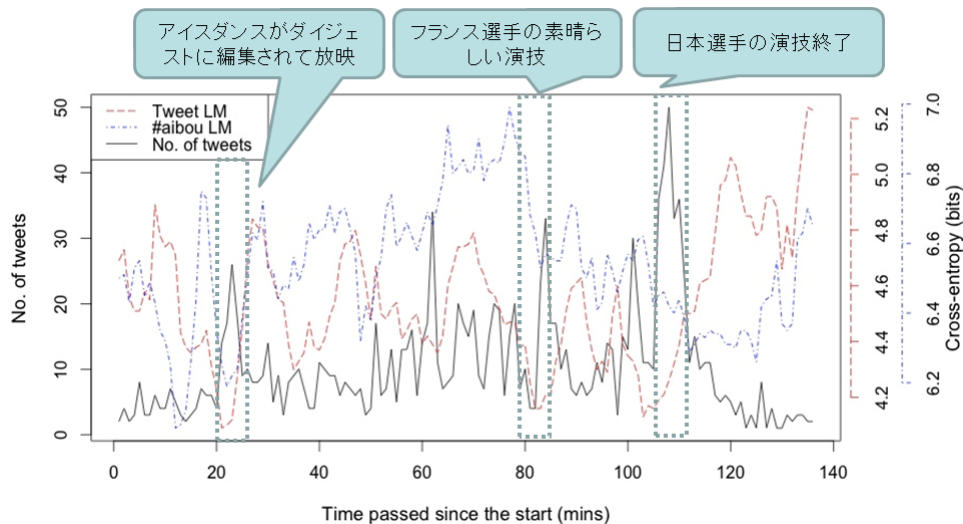


図 5 #figureskate データに対し、ツイート言語モデルと#aibou 言語モデルをそれぞれ用いてクロスエントロピーを算出した結果.

Fig. 5 The occurrence of #figureskate data shown together with the cross-entropy between the general tweet language model and the #aibou language model.

ロピーが低くなった盛り上がり時間帯の中で、低くならなかった時間帯がある。図 4 のイベントの始まった瞬間にツイート言語モデルではクロスエントロピーが低下したが、#aibou 言語モデルを用いた際には低下しなかった。このことから、#aibou 言語モデルはジャンルの異なるハッシュタグに対して盛り上がり時間帯を捉えることに向いていないことがわかる。

そのため、各ハッシュタグにおいて盛り上がり時間帯とその他の時間帯よりそれぞれ 150 ツイート、合計 300 ツイートをサンプリングした。サンプリングしたツイート集合に対し、ツイート言語モデルを用いて計算したクロスエントロピーを表 3 に示す。表 3 の 4 列目に示した値は盛り上がり時間帯とその他の時間帯との二つのツイート集合

に対して、Wilcoxon の順位和検定に算出した  $p$  値である。また表 3 に示した 6 つのハッシュタグに対して、盛り上がり時間帯とその他の時間帯におけるクロスエントロピーの差に対し Wilcoxon の符号順位検定を適用した結果、統計的に有意差が存在した ( $p < 0.02$ )。n-gram 以外の特徴量も影響していると推測し、5 章でより詳細に解析した。

## 5. n-gram 以外の特徴の検討

n-gram 以外の特徴を考慮するために、教師ありデータを用いて入力ツイートが盛り上がり時間帯におけるものか否かの二値分類を行うモデルを構築した。採用した学習モデルは SVM と Random Forest であり、実装にはそれぞれ LIBSVM パッケージ [4] と RandomForest パッケージ [13]

表 3 各ハッシュタグの盛り上がり時間帯とその他の時間帯よりサンプリングしたツイート集合とツイート言語モデルを用いて算出したクロスエントロピー。

Table 3 Cross-entropy of sampled set of tweets from spikes and non-spikes computed using the general tweet language model.

ハッシュタグ	H(盛り上がり時間帯)	H(その他の時間帯)	H(その他の時間帯) - H(盛り上がり時間帯)	p 値
#aibou	5.45	5.58	0.13	$p < 0.3$
#hanshin	5.54	6.15	0.61	$p < 0.0004$
#ACV	5.44	5.61	0.17	$p < 0.06$
#agqr	4.22	5.11	0.89	$p < 10^{-9}$
#figureskate	4.18	4.78	0.60	$p < 10^{-11}$
#momoclo	4.73	5.17	0.44	$p < 0.00004$

を用いた。

### 5.1 使用した素性

使用した素性は以下の通りである。

- (1) 漢字の数
- (2) ひらがなの数
- (3) カタカナの数
- (4) ツイート中の合計文字数
- (5) 各文字ユニグラム数
- (6) 同一文字の3文字繰り返し回数

Brody ら [2] の研究に基づき、盛り上がり時間帯において極性を持つ単語が現れやすいと推測したため、3文字繰り返し回数を素性として入れた。また、学習データが300件と少ないので、 $n \geq 2$  の n-gram 素性は扱わなかった。

### 5.2 結果と考察

二値分類において正しく盛り上がり時間帯のツイートとして分類した数と誤って盛り上がり時間帯のツイートとして分類した数を考慮するため評価指標として Accuracy (Acc.) を用いた。Accuracy の算出式を以下に示す：

$$\text{Acc.} = \frac{\text{正しく分類した数} + \text{誤って分類した数}}{\text{Total No. of tweets}}$$

表 4 に SVM によるツイートの分類結果を示す。この実験では1ハッシュタグからサンプリングされたデータを学習データとして用い、残りの5ハッシュタグのデータをテストデータとした。

表 5 に各ハッシュタグに対して Random Forest を適用した結果算出された重要な特徴量の上位5件を示す。各特徴量に対して、Wilcoxon の順位和検定を適用した結果、統計的有意差を示した ( $p < 0.01$ )。Random Forest では n-gram 言語モデルでは捉えられなかった特徴量である漢字数が重要な特徴量であると捉えられている。5ハッシュタグにおいて漢字数が上位3位以内にある。なお#agqrにおいて漢字数は上位5位に含まれていないが、7番目であった。また、文字ユニグラムは各ハッシュタグに依存していることがわかる。例を挙げると#ACVにおいて上位1位、2位にある文字ユニグラムはサーバが落ちたことによって

ツイート数が急激に増加したことに原因があるため、「落」と「ち」が上位に位置している。

表 4 より、#aibou データにおいてクロスエントロピーが低いにも関わらず、SVM による分類では高い性能を記録している。このことからクロスエントロピーのみでは盛り上がっている時間帯におけるツイートの性質を捉えきれていない、もしくはデータセットを構築する際のサンプリング手法においてノイズが含まれていると考えられる。

表 5 から盛り上がっている時間帯のツイートの性質として漢字数が少ないことが明らかになった。工藤ら [22] が述べているように、ユーザはツイートする時、漢字を入力するために Input Method Editor (IME) を利用しているため、ひらがなを入力する際より打鍵コストが多くかかる。そのため、漢字数が少ないという特徴は打鍵コストが関係しているのではないかと考えられる。

また、漢字の種類数はひらがな、カタカナと比較すると種類が多いことから一般的に高い情報量を保持している。しかし Web 上に用いられる言葉は特定の単語においてはひらがなが用いられやすいことが知られている [22]。例えば、「なく」という単語は「泣く」「無く」「鳴く」等の同音異義語が存在するが、即時性が要求される場面である盛り上がり時間帯では漢字が用いられにくい。以上のことから、盛り上がり時間帯におけるツイートは Twitter 全体からサンプリングしたツイートデータをもとに構築したツイート言語モデルに比べ、漢字数が少なく、クロスエントロピーが低い、と推測できる。

## 6. 結論

本稿では盛り上がっている時間帯のツイートはツイート言語モデルによって算出したクロスエントロピーが低いことを示した。SVM による分類結果より分類性能はハッシュタグに依存することがわかる。また、Random Forest において重要な特徴量の解析を行った結果、盛り上がっている時間帯のツイートにおいて漢字数が少ないことも示した。このことから盛り上がり時間帯におけるツイートは打鍵コストが関係しているのではないかと推測できる。

今後の課題として、1) クロスエントロピーが低いこと、



- of the Association for Computational Linguistics, Los Angeles, California, Association for Computational Linguistics, pp. 181–189 (online), available from (<http://www.aclweb.org/anthology/N10-1021>) (2010).
- [16] Sakaki, T., Okazaki, M. and Matsuo, Y.: Earthquake shakes Twitter users: Real-time event detection by social sensors, *Proc. of the 19th International Conference on World Wide Web*, ACM, pp. 851–860 (2010).
- [17] Sasano, R., Kurohashi, S. and Okumura, M.: A simple approach to unknown word processing in Japanese morphological analysis, *Proc. of the Sixth International Joint Conference on Natural Language Processing*, Nagoya, Japan, Asian Federation of Natural Language Processing, pp. 162–170 (online), available from (<http://www.aclweb.org/anthology/I13-1019>) (2013).
- [18] Shen, C., Liu, F., Weng, F. and Li, T.: A participant-based approach for event summarization using Twitter streams, *Proc. of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Atlanta, Georgia, Association for Computational Linguistics, pp. 1152–1162 (online), available from (<http://www.aclweb.org/anthology/N13-1135>) (2013).
- [19] Thelwall, M., Buckley, K. and Paltoglou, G.: Sentiment in Twitter events, *Journal of the American Society for Information Science and Technology*, Vol. 62, No. 2, pp. 406–418 (2011).
- [20] Zhao, W. X., Jiang, J., Weng, J., He, J., Lim, E.-P., Yan, H. and Li, X.: Comparing Twitter and traditional media using topic models, *Proc. of the 33rd European Conference on Advances in Information Retrieval*, Berlin, Heidelberg, Springer-Verlag, pp. 338–349 (online), available from (<http://dl.acm.org/citation.cfm?id=1996889.1996934>) (2011).
- [21] 齊藤いつみ, 貞光九月, 浅野久子, 松尾義博: 正規-崩れ表記のアライメントに基づく表記崩れパタンの抽出と形態素解析への導入, 情報処理学会研究報告 NL214, pp. 1–9 (2013).
- [22] 工藤 拓, 市川 宙, Talbot, D., 賀沢秀人: Web 上のひらがな交じり文に頑健な形態素解析, 言語処理学会第 18 回年次大会, pp. 1272–1275 (2012).