

タイトル	現代日本語書き言葉均衡コーパスコアデータにおけるオノマトペ出現実態に基づくオノマトペ自動抽出手法
著者	内田, ゆず; Uchida, Yuzu
引用	工学研究 : 北海学園大学大学院工学研究科紀要(17): 15-20
発行日	2017-09-30

研究論文

現代日本語書き言葉均衡コーパスコアデータにおける
オノマトペ出現実態に基づくオノマトペ自動抽出手法

内 田 ゆ ず*

Onomatopoeia Extraction Method
Based on Usage of Onomatopoeias in BCCWJ Core Data

Yuzu Uchida*

概 要

近年、オノマトペに関する研究が発展している。しかし、現代のオノマトペ使用実態を反映した体系的なオノマトペ辞書は整備されていない。筆者らは、現実に使用されているオノマトペとその用例文を収集し、大規模な辞書アプリケーションを構築しようとしている。本稿では、現代日本語書き言葉均衡コーパスのコアデータを対象として、オノマトペの出現実態や品詞情報を分析した結果について報告する。さらに、分析により得られた知見に基づいて、オノマトペの自動抽出手法を提案する。

1. はじめに

オノマトペ（擬態語、擬音語の総称）は、自然界の音や事物・動作の様態を表す語群で、日本語の語彙に豊富に存在している。近年、オノマトペを様々な分野で利活用することを目指した研究が進められている¹⁾。

オノマトペには多様な語義をもつという特徴がある。例えば日本語オノマトペ辞典²⁾の「ごろごろ」の項目には6つの語義が掲載されている（「雷の響く音」「猫がのどを鳴らす音」等の擬音の語義と「無造作に転がっているさま」「仕事をせずに無駄に暮らしているさま」等の擬態の語義）。高丸らによる地方議会会議録コーパスにおける「ごろごろ」を含む文の分析では、辞典中の語義に加えて「たくさんある（いる）さま」、「変わりゆくさま」などの語義が見られた。このように、1つのオノマトペが擬音、擬態の語義を持つだけでなく、さらにそれらから派生した語義や新たな語義が追加されることが報告されている³⁾。また、語義が類似したオノマトペが多数あるという特徴もある。例えば「ごろごろ」に対して、「ころころ」「ごろんごろん」「ごろっ」は類似の語義をもつものの、

それらが表現する様子やニュアンス、修飾できる語はそれぞれやや異なると予想される。

これらのことは日本語母語話者にとっては直感的に理解可能であるが、日本語学習者にとっては理解が容易ではない。また、自然言語処理などの工学分野においてオノマトペを利用する場合にも、語義の曖昧さが障壁となる。日本語非母語話者がオノマトペを適切に使用するためには、あるオノマトペがどのような場面で使用可能であるかという実例を示すことが重要であるし、自然言語処理分野のタスクにおいても、前方および後方の文脈に基づいて、適切なオノマトペを選択する必要があると考えられる。

そこで、筆者らはオノマトペの実際の用例を対象として「オノマトペの語義」と「オノマトペと共起する語（コロケーション）」に着目した研究を進めている。オノマトペを含む用例文から、例えば「ごろごろ」+「寝る」、「ごろごろ」+「転がる」という係り先のコロケーションや、「石が」+「ごろごろ」、「雷が」+「ごろごろ」、「のどを」+「ごろごろ」という係り元のコロケーションを抽出し、そのオノマトペが使用できる文脈を明らかにする。人間がオノマトペを学習する際には、さ

* 北海学園大学大学院工学研究科電子情報生命工学専攻

Graduate School of Engineering (Electronics, Information, and Life Science Eng.), Hokkai-Gakuen University

らに各コロケーションの具体的な例文を提示することで、語義を計り知ることが可能であろう。また、「Aが」+「ごろごろ」+「転がる」と、「Bが」+「ころころ」+「転がる」という共起を考えたときに、「ごろごろ」と共起する単語集合Aと、「ころころ」と共起する単語集合Bの差異を見れば、2つのオノマトペの意味の違いを理解することにつながると考えられる。このような観点から、本研究では現代のオノマトペの最新の用法を提示できるウェブ上の実例に基づく辞書（オノマトペ実例辞書）の構築を目指している。オノマトペ実例辞書構築のためには、現代日本語における最新の用例が多数必要である。そこで、ウェブ上の文書からオノマトペを抽出し、「オノマトペ用例データベース」を構築する。オノマトペは文字長の短いひらがな／カタカナの文字列であり、特殊拍（促音・撥音・長音）の挿入により変形が可能であるため、文書中からオノマトペを正確に抽出することは難しい。ブログや議会会議録からオノマトペを自動抽出する手法が提案されている⁴⁾⁵⁾⁶⁾が、更なる検討が必要である。

次に、「オノマトペ用例データベース」内の文に対して、係り受け解析や共起する単語の纏め上げを行うことで、オノマトペ実例辞書に必要なコロケーションデータを得る。大規模言語資源とコロケーションに関する研究はこれまでも行われている。田野村は、ウェブコーパスから得られるコロケーション情報からのコロケーション辞典作成の手法について幾つかの具体例を元に考察している⁷⁾。部は、「しんみり」「しみじみ」の2語を対象に新聞コーパスにおけるコロケーション（共起する動詞）を調査し、アンケート調査によって得た人間が想起する係り先の動詞と比較している⁸⁾。

本稿では、オノマトペコロケーション抽出の出発点として、「現代日本語書き言葉均衡コーパス(BCCWJ)」のコアデータに含まれる全てのオノマトペの表層形態を分析する(3章)。さらに、この結果に基づき、品詞情報を利用してオノマトペの抽出を行い(4章)、抽出手法の拡張を試みる(5章)。最後に、コアデータから得られるコロケーションの例について触れつつ結論を述べる(6章)。

2. 対象データ

本研究で使用するデータについて説明する。

2.1 コーパス

本研究で分析対象とするコーパスは、大学共同利用機関法人人間文化研究機構国立国語研究所と文部科学省科学研究費特定領域研究「日本語コーパス」プロジェクトが共同で開発した『現代日本語書き言葉均衡コーパス』⁹⁾(Balanced Corpus of Contemporary Written Japanese, 以降BCCWJ)である。BCCWJには、現代の日本語の書き言葉の全体像を把握できるように集められたサンプルが書籍全般、雑誌全般、新聞、白書、ブログ、ネット掲示板、教科書、法律などのジャンルにまたがって約1億430万語収録されている。

なお、BCCWJには人手で形態素解析結果を修正したサブセットであるコアデータが含まれている。コアデータは約9万短単位のデータである。

2.2 オノマトペ辞典

ある単語がオノマトペであるかを判断する際に、日本語オノマトペ辞典²⁾を基準として用いる。この辞典には古事記などの古典から現代に至るまでのオノマトペが掲載されており、見出し語の数は4,564語となっている。

この辞典には2種類の索引がある。一つ目の「意味分類別さくいん」は、辞典に収録されている見出し語のうち、延べ2,470語(異なり1,751語)を採り上げ、自然・人間・事物に三分類し、それぞれに簡略な解説を付してあるものである。一般性の高い語が厳選されたオノマトペ集合と考えられる。二つ目の「五十音順さくいん」には、辞典の本編ならびにコラム、付録(漢語オノマトペ、鳴き声オノマトペ)に収録されている全4,506語が掲載されている。漢語オノマトペを含む表現(例:焔焔に滅せずんば炎炎を若何せん)や、オノマトペではないがコラムで言及されている語(例:あいまい)も対象であるため、語数は多いがオノマトペとして不適切なものも含まれている。

3. コアデータの全オノマトペ分析

筆者らは本研究に先立ち、BCCWJに出現するオノマトペの傾向を概観するため、意味分類別さくいんに掲載されたオノマトペ(1,751語)と完全一致する短単位形態素をコアデータからすべて抽出し、分析を行った。その結果、抽出された短

あつあつ、あつさり、いちやいちや、がが、かくかく、
 かつちり、がつつり、がらっ、がらん、ぎざぎざ、ぎっ
 くり、きつちり、きりっ、きりり、ぐいっ、くつきり、く
 ったり、くるくる、ぐるぐる、くるっ、ぐるっ、ぐるり、
 くんくん、ぐんなり、ごうごう、こじんまり、こぢんま
 り、こっそり、こてんばん、ささっ、しっくり、じっと、
 しゃなり、しゅわしゅわ、ずしり、すっかり、すつきり、
 ずっしり、ずっと、すっぼり、すばすば、すべすべ、す
 ぼっ、するっ、すれすれ、そっくり、そろり、だらり、ち
 ぐはぐ、ちゃんと、ちゅんちゅん、ちよい、ちよくちよ
 く、ちよこまか、ちよんちよん、つるっ、でこでこ、て
 つきり、でれでれ、でん、どきん、とことん、どっかり、
 どっしり、とんかち、どんびしゃ、によるによる、ばか
 っ、ばっくり、ばったり、ばっちり、ばばば、ばん、びし
 ばし、ひしひし、ひっそり、びび、ひよっこり、びよん
 びよん、おすん、ぶちぶち、ふつつつ、ふらっ、ふらり、
 ペこり、ペしゃペしゃ、ペしゃんこ、べちやり、ぼうぼ
 う、ぼか、ぼちっ、ぼちぼち、ほっこり、ほっぼ、ほつり
 ほつり、まったり、まんまん、むちゃくちゃ、もちもち、
 もっちり、もんもん

図1 意味さくいん：掲載なし／
五十音さくいん：掲載ありのオノマトベ

単位は5,133個であり、そのうち1,370個がオノマトベであった(異なり数:392語)。つまり、オノマトベと字面が一致する短単位のうち、73.3%はオノマトベではないことになる。この点について、オノマトベの文字数別に集計すると顕著な傾向が見られる。2～3文字のオノマトベと一致する短単位の93.0%、4文字以上のオノマトベと一致する短単位の5.9%がオノマトベではなかった。ここから、2～3文字の短いオノマトベを抽出するためには単に表層を手がかりにするのではなく、品詞等の情報が必要であることが明らかになった。したがって、本章ではコアデータ中の全てのオノマトベを抽出し、それらの品詞を分析する。

コアデータ中の2文字以上のひらがな・カタカナからなる短単位形態素を全て抽出し、それらがオノマトベであるかを人手で判断する。この分析によって、コアデータ中の全てのオノマトベ(つまり、正解データ)を得ることを意図している。

分析の結果、198,829個の短単位が抽出され、そのうち2,048個がオノマトベであると判断された。意味分類別さくいんに掲載されていないオノマトベは182語、延べ678回出現している。

意味分類別さくいんには掲載されていないが五十音順さくいんに掲載されているものは図1に示す101語である。これらのオノマトベは、五十音順さくいんを導入することで抽出が可能になる。

うがうが、うんと、かあん、がたがたがた、がちゃ、か
 っかっかっ、かっつ、きちつと、きちんと、ぎゃあぎゃ
 あ、きゃつきゃ、ぎゅ、ぎよつと、ぐらぐらぐら、くり
 くりっ、ぐるぐるぐるっ、ぐんと、こつ、ごふつ、こり
 こりっ、さささささっ、さっさと、しんと、じんと、す
 うっ、ずうつと、ずず、すつかり、ずらずらっ、せつせ
 と、そつと、ぞつと、たたたつた、ちびり、ちよいと、
 ちよこつと、ちら、つるんつるん、てれん、てんかん、
 とつとと、どよどよ、とんとんとん、のうのう、はたと、
 ばっさり、ばっさ、ばっちし、はつと、ばぼっ、ぼらぼ
 らっ、ばんっ、びいい、ひいひい、びか、びくと、びぼび
 ぼ、ひよつと、ふうふう、ぶふおつ、ぶぶっ、ぶらつと、
 ふらふらっ、ぶろろろ、ぶんすか、ペこ、べたりんちよ、
 べろりっ、ぼうつと、ぼそ、ほたり、ぼち、ほつと、ぼに
 よ、ぼぼん、ほわり、むっちゃ、めちゃ、めっちゃ、めっ
 ちゃめちゃ、よよと

図2 意味さくいん：掲載なし／
五十音さくいん：掲載なしのオノマトベ

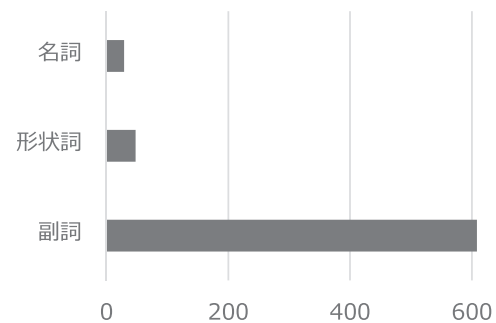


図3 オノマトベの品詞

意味分類別さくいんにも五十音順さくいんにも掲載されていない語は図2に示す81語である。「きちんと」や「くりくりっ」はそれぞれ索引に掲載された「きちん」、「くりくり」に助詞「と」、促音「っ」を付与することで対応できる。このように、一部のオノマトベは単純なルールで抽出が可能になる。一方、「ごふっ」や「ぶんすか」などは比較的新しい表現だと考えられ、このような新出オノマトベを抽出する手法の確立が求められる。

図3にオノマトベであると判断された短単位の品詞を示す。すべての短単位が副詞、形状詞、名詞のいずれかに分類され、88.8%は副詞である。品詞を抽出の条件に加えることで、短いオノマトベの抽出精度を向上させることが期待できる。

4. 品詞情報を利用したコアデータからのオノマトベ抽出

3. の結果に基づき、品詞情報を利用したオノ

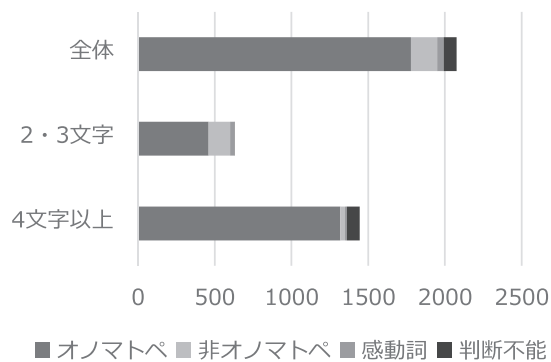


図4 コアデータからのオノマトペ抽出結果

マトペの抽出実験を行う。コアデータに MeCab¹⁰ (Unidic 辞書) で形態素解析を施し、五十音順さくいんに掲載されているオノマトペと字面が一致し、かつ副詞あるいは形状詞になった短単位を人手で分析する。

図4に抽出結果を示す。2,076個の短単位が抽出され、1,778個(85.6%)がオノマトペであった。品詞情報を用いることでオノマトペを高い精度で抽出できることが明らかになった。しかし、3.では考慮していなかった形態素解析誤りや対象オノマトペの拡充に起因するエラーが発生した。

以下に人手で非オノマトペと判断された例をエラーの原因ごとに示す。(下線部が該当箇所)

- ① 2文字/長音/カタカナ (形態素解析誤り)
 - ・ 育ち盛りの高校生、こーゆー添加物のこと…
 - ・ …おともだちがサッかーのしあいがありました。
 - ・ 一番目立っていたグレートデン。
 - ・ どーでもイイ。
 - ・ …限定販売する「ビープラスDT」(16万円)だ。
- ② 助詞との接続 (形態素解析誤り)
 - ・ 挽き出すときに、目がちゃっとひっかかるわけですわ。
 - ・ 病気のペット (たとえばワンちゃんとしましよう)は…
- ③ コラム掲載語
 - ・ 責任もあいまいだった。
 - ・ こわごわ組んだローンだけど…
 - ・ わたしは、みにくい姿の魔物がすきだ。
 - ・ フルに使いこなすには取説が必要かも。

④ 同音異義語

- ・ 私にはたった一つだけ望みがあった。
- ・ …おうおうにして東洋趣味に走るのよね。
- ・ 「かくかくしかじか？」で…
- ・ これが一般人のごくごく健全な感覚でしょう。
- ・ 二十年も放置され、とうとう空家が一千戸に達した。
- ・ 若い人たちの話をよくよく聞いてみると…

判断不能とされたのは、「しばしば」、「だんだん」、「まだまだ」、「みすみす」、「みるみる」など、一般の副詞として認識されつつあるオノマトペである。

この実験の結果から、本手法の改善には、五十音順さくいんから一部の語を除くことや、同音異義語の問題を回避するためにストップワード(オノマトペと品詞の組)を設けることが有効だと考えられる。

5. BCCWJのブログデータからのオノマトペ抽出

4.で述べた改善策を導入したオノマトペの抽出手法を提案する。解析誤りが特に起こりやすいカジュアルな文体での本手法のオノマトペ抽出精度を確認するため、BCCWJに含まれるYahoo!ブログのデータを対象として抽出実験を行う。具体的なアルゴリズムを図5に示す。

提案手法によって、49,492個の短単位がオノマトペとして抽出された。これまでの分析で、2～3文字のオノマトペの抽出精度が特に低いということが明らかになっている。したがって、ここでは抽出された短単位のうち、2～3文字のものを全て人手で確認し、オノマトペであるかを判断する。

2～3文字の短単位は14,706個抽出され、9,749個(66.3%)がオノマトペであった。2文字の短単位のみでは37.8%、3文字のみでは86.0%の精度である。表1に各ルールにおける抽出結果を示す。

2文字のオノマトペの抽出エラーが全体の抽出精度に悪影響を及ぼしていることがわかる。中でも、2文字の短単位にルールbを適用すると、非オノマトペをオノマトペとして抽出するエラーが多い。これは、形態素解析誤りによって別の単語

表1 各ルールによるプログデータからのオノマトベ抽出結果

	ルール a		ルール b		ルール c		ルール d		ルール e		ルール f		計	
	2文字	3文字	2文字	3文字	2文字	3文字	2文字	3文字	2文字	3文字	2文字	3文字	2文字	3文字
オノマトベ	0	13	1689	5312	583	21	0	40	0	1921	5	165	2277	7472
非オノマトベ	0	19	2356	942	873	72	0	8	0	122	152	24	3381	1187
感動詞	0	0	0	1	14	0	0	0	0	0	0	0	14	1
判断不能	0	0	296	11	10	0	0	0	0	16	41	0	347	27
計	0	32	4341	6266	1480	93	0	48	0	2059	198	189	6019	8687

- I. 3種類のリストを作成する
- ・オノマトベリスト：五十音順さくいんから不適切な語を除いたリスト
 - ・品詞例外リスト：これまでの分析で明らかになった、副詞・形状詞以外に分類されるオノマトベとその品詞をペアのリスト
 - ・ストップワードリスト：これまでの分析で明らかになった、オノマトベとの同音異義語のリスト
- II. MeCab (Unidic 辞書) で形態素解析を行う
- III. オノマトベリスト中の語と字面が一致する短単位を抽出する
- IV. III で抽出された短単位のうち、以下の条件を満たすものをそれぞれオノマトベと判断する (抽出ルール)
- 品詞が副詞、形状詞以外で、品詞例外オノマトベリストに存在する
 - 品詞が副詞か形状詞で、オノマトベリスト中の語と完全一致し、ストップワードリストに存在しない
 - 品詞が副詞か形状詞で、オノマトベリスト中の語から最終促音を削除したものと一致し、ストップワードリストに存在しない
 - 品詞が副詞か形状詞で、オノマトベリスト中の語に最終促音を付加したものと一致する
 - 品詞が副詞か形状詞で、オノマトベリスト中の語に助詞「と」を付加したものと一致する
 - 品詞が副詞か形状詞で、長音母音を長音記号に変換、あるいは繰り返しの縮約を行うとオノマトベリスト中の語と一致する

図5 オノマトベ抽出アルゴリズム

の一部や長いオノマトベの一部が切り出されることが主な原因である。具体例を以下に挙げる。例中の下線部が該当箇所、すべて形態素解析によって副詞と判断されている。

- ・あなたがた自身も、あらゆる行いにおいて…
- ・よし、少しベンキョーすっかな…まず、初めに…
- ・ぶりんっとしたやつね。

オノマトベによっては、定型句のような表現で

しか使われないものもある。形態素解析誤りは避けられないので、個別のオノマトベに抽出ルールをカスタマイズするなど、品詞情報だけに頼らない抽出手法を検討したい。

また、ルール c~f は、辞典に掲載されていない新しいオノマトベや、既存のオノマトベが変形したものを抽出する役割を果たすが、過剰に適用されることでエラーの原因にもなる。具体例を以下に挙げる。

- ・どなたかコツ教えて下さい。

(ルール c 適用：「こつっ」の促音削除)

- ・…するまでほっとうと思って今になりました。

(ルール e 適用：「ほっ」に「と」を付加)

- ・友達は私といるだけで巻き添えくうし…

(ルール f 適用：「くー」に変換)

- ・気付けばあっと言う間の12月。

(ルール f 適用：「ばっ」に変換)

- ・しかも土曜日にウチに来て←パパン居なかったから

(ルール f 適用：「パン」に縮約)

これらのルールの精度向上には、係り受け関係などの利用や、ルールの適用範囲を限定する工夫が必要だろう。

6. おわりに

本稿では、オノマトベ実例辞書の構築を目指し、BCCWJ を対象としたオノマトベの抽出及び分析を行った。オノマトベ—特に短いオノマトベ—の自動抽出には多くの課題が残されているものの、豊富なオノマトベの実例を得た。

それらの実例の中には、いくつかの興味深いコロケーションが見受けられる。たとえば、「ぐんぐん」は「伸びる」や「大きくなる」などの成長に関わる動詞に係ることが多い、「くるくる」や「ぐるぐる」は回転に関わる動詞に係る点で共通して

いるが、「ぐるぐる」は「さまよう」などの動詞に係ることもある、「しゅわしゅわ」は炭酸入りの飲料水とともに用いられることが多いなどである。

現状では、オノマトベと共起する表現を纏め上げてコロケーションを得るには実例が不足している。今後、ウェブ上のデータを収集・整理した、地方議会会議録コーパス、ブログコーパスを利用し、大規模なコロケーションデータベースの構築を行う予定である。

謝辞

本研究は北海学園学術研究助成および科研費(No.26370498)の助成を受けたものである。また、本研究は宇都宮共和大学の高丸圭一准教授、福岡大学の乙武北斗助教、小樽商科大学の木村泰知准教授と共同で推進された。

参考文献

- 1) 小松孝徳：論文特集「オノマトベの利活用」にあたって、人工知能学会誌 30(1), p.134, 2015.
- 2) 小野正弘編：日本語オノマトベ辞典, 小学館, 2007.
- 3) 高丸圭一, 内田ゆず, 乙武北斗, 木村泰知：地方議会会議録コーパスにおけるオノマトベ出現傾向と語義の分析一, 人工知能学会論文誌, 30(1), pp.306-318, 2015.
- 4) 内田ゆず, 荒木健治, 米山淳：ブログ記事からのオノマトベ用例文の自動抽出手法, Journal of Japan Society for Fuzzy Theory and Intelligent Informatics, 24(3), pp.811-820, 2012.
- 5) 木村泰知, 渋谷英潔, 内田ゆず, 乙武北斗, 高丸圭一, 森辰則：地方議会会議録におけるオノマトベの自動抽出手法の提案, 第30回ファジィシステムシンポジウム講演論文集, pp.638-641, 2014.
- 6) 池田祐一, 阪本浩太郎, 渋谷英潔, 森辰則：国際音声記号を素性とした3文字以下の未知のオノマトベ自動抽出手法の提案, 言語処理学会第21回年次大会論文集, P1-12, 2015.
- 7) 田野村忠温：日本語コーパスとコロケーション一辞書記述への応用の可能性一コーパスからのコロケーション情報抽出一分析手法の検討とコロケーション辞典項目の試作, 阪大日本語研究, 21, pp.21-41, 2009.
- 8) 部楓：コーパスを利用した類義語のコロケーション分析一擬態語「しんみり, しみじみ」と動詞の共起から一, ことばの科学, 19, pp.129-140, 2006.
- 9) 山崎誠編：『書き言葉コーパス一設計と構築一』講座日本語コーパス2, 朝倉書店, 2014.
- 10) Taku Kudo, Kaoru Yamamoto, Yuji Matsumoto: Applying Conditional Random Fields to Japanese Morphological Analysis, Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP-2004), pp.230-237, 2004.