

令和 2 年 6 月 19 日現在

機関番号：62618

研究種目：基盤研究(B)（一般）

研究期間：2015～2019

課題番号：15H03210

研究課題名（和文）統語・意味解析情報タグ付きコーパス開発用アノテーション研究：複文を中心に

研究課題名（英文）Research on annotation for the development of a parsed corpus of Japanese with a special focus on complex sentences

研究代表者

PARDISHI P.V. (Pardeshi, Prashant)

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・理論・対照研究領域・教授

研究者番号：00374984

交付決定額（研究期間全体）：（直接経費） 12,800,000円

研究成果の概要（和文）：本研究は現代日本語の特徴の一つである関係節および従属節を中心とする複文について言語学的情報を検索、抽出するために必要なアノテーション方法を研究し、それに基づいてタグ付け作業を行い、複文に関するより高度の環境整備を整えることを目標とした。この作業は国立国語研究所の「統語・意味解析コーパスの開発と言語研究」プロジェクトと連携をしながら進め、2020年3月時点で、40,831万文規模のコーパス(560,098語)を検索ツールとともに公開した。

研究成果の学術的意義や社会的意義

本研究の成果として約4万文(56万語)規模の統語・意味解析付きコーパスが開発され、公開された。このコーパスでは現代日本語のテキストに対し文の統語・意味解析情報を付与し、複数の検索ツールが用意されており、多様な日本語の機能語や句構造、節の諸類型および複雑な構文を大量の言語データから検索・抽出して研究に活用することが可能である。また、初心用の検索ツール完備により、大学などで日本語の統語論の教育にも利用できる。さらに、このコーパスでは日英両言語表記で公開されており、日本語表記に精通してない海外の研究者も利用できるため、日本語研究の国際化に貢献できる。

研究成果の概要（英文）：The goal of this research project was to build a large-scale parsed corpus (treebank) of modern Japanese that would enable the users to search and retrieve complex sentences involving relative clauses and subordinate clauses, which are peculiar characteristic of Japanese, through the development of annotation method and corpus search tools. This task was jointly carried out with the collaborative research project “Development of and Linguistic Research with a Parsed Corpus of Japanese” at NINJAL. A corpus of 40,831 sentences (560,098 words) was build and is made available for free access at the following site: <http://npcmj.ninjal.ac.jp/>. The corpus can be searched with the following search tools: <http://npcmj.ninjal.ac.jp/explorer/> [http://npcmj.ninjal.ac.jp/interfaces/index\\_en.html](http://npcmj.ninjal.ac.jp/interfaces/index_en.html)

研究分野：コーパス言語学

キーワード：関係節 複文 アノテーション 統語・意味解析付きコーパス コーパス研究

## 1. 研究開始当初の背景

本研究の代表者は第二期中期計画において『現代日本語書き言葉均衡コーパス (BCCWJ)』のデータを利用して、内容語(動詞、名詞、形容詞、副詞)の振る舞い(主に共起関係)を検索・抽出するシステム NLB を開発した(<http://nlb.ninjal.ac.jp/>)。この検索システムでは、特に高度なコンピュータの知識がないユーザーも、簡単かつ瞬時に研究に必要な情報を検索・抽出・ダウンロードでき、ユーザー数は数万人となっている。この間、国内外の研究者・大学院生からコーパスに基づく機能語、句、節、複文といった様々なレベルでの構造体(constructions)を検索・抽出できるシステムを構築してほしいとの要望を多数受けた。また、学会などでも研究者らからその要請を受けた。

英語に関しては、90年代初頭から Pennsylvania (以下、Penn) 大学で統語解析情報タグ付きコーパスが開発され、20年かけて意味情報も付加したコーパスが複数完成し、コーパスに基づく機能語、句、節、複文のような様々なレベルでの構造体を大量のデータから検索・抽出して研究を行うことが可能となり、目覚ましい成果をあげている。

英語の成功を受けて、世界各国では、言語研究に資するコーパスが開発されている。共通なフォーマットで構築される Penn 方式のコーパス間での対照研究も活発に行われている。一方、日本語に関しては、言語研究を目的とする汎用的統語・意味解析情報タグ付きコーパスは未だに存在せず、コーパスに基づく日本語の諸構造体の研究は英語などと比べて、立ち遅れている。

## 2. 研究の目的

上記の問題意識を踏まえ、本研究は現代日本語における複文の一つである連体節を検索・抽出するために必要なアノテーション方法を研究し、それに基づいて約5万文(約90万語)に対してタグ付け作業を行い、連体節を簡単に検索・検出できる環境整備を整えることを目標とした。

## 3. 研究の方法

統語解析情報タグ付きコーパス開発のためには日本語の諸構造体(constructions)を検索・検出できるためのタグ付け(アノテーション)作業が不可欠である。この基礎作業を行うためには日本語文法の専門的知識、特に統語的な側面の専門知識が必要であり、さらに、大量データに対してタグ付け作業を(半)自動的に行う必要があるため、形式文法、コンピューター言語学および自然言語処理といった言語関連諸分野の共同研究・共同作業を行った。

具体的には、本研究は現代日本語の特徴の一つである関係節および従属節を中心とする複文について言語学的情報を検索、抽出するために必要なアノテーション方法を研究し、それに基づいて約4万文(約55万語)に対してタグ付け作業を行い、複文に関するより高度の環境整備を整えた。なお、この作業は国立国語研究所の「統語・意味解析コーパスの開発と言語研究」プロジェクトと連携をしながら進め、国立国語研究所のホームページでコーパスを一般公開した(詳細は以下を参照)。

## 4. 研究成果

統語解析情報タグ付きコーパス開発とそれを利用して研究を行い、以下のような、成果がありました。

### ● コーパスおよび検索ツールの構築と公開

NPCMJ (NINJAL Parsed Corpus of Modern Japanese) コーパス構築：研究期間を通して、毎年アノテーション作業を行い、では2020年3月現在、約4万文(ツリー)、語数にして約56万語のデータが収録されているコーパスを構築し、国立国語研究所のウェブサイトを通じて(<http://npcmj.ninjal.ac.jp/>)無償公開した。

データはすべて、漢字仮名混じりとローマ字の両方で表記されており、全データ、あるいは下で紹介する検索ツールを使った検索結果を自分のコンピューターにダウンロードすることが可能である。

NPCMJ コーパス検索用に専用の検索ツールの開発も行い、初中級者向けの NPCMJ Explorer および中上級者向けの NPCMJ Search を公開した。

NPCMJ Explorer は、単純な文字列検索と、特定の文法項目に関する用例の検索という2種類の検索手段を備えたツールであり、文法項目として、益岡隆志・田窪行則(1992)『基礎日本語文法』(くろしお出版)で解説されている136の項目から73項目を取り上げ、各項目のラベルをクリックするだけで、利用者が複雑な検索式を作成しなくても、用例が表示される仕組みになっている。ジャンル(表1における「出典」)ごとの頻度を調べるためのジャンル指定機能もある。

NPCMJ Search では、品詞や句・節に与えられたラベルや、ラベルとラベルの関係を検索式

によって指定しながら用例の検索が可能である。どちらのツールでも、検索結果を一覧表示させたり、木構造で表示させたりすることができる。本コーパス・検索ツールの開発によって複文に関するより高度の環境整備を整えることができた。

- 雑誌論文・学会発表

本研究の成果を国内外の学会や研究雑誌などを通じて公開した。

## 5. 主な発表論文等

〔雑誌論文〕 計11件（うち査読付論文 0件 / うち国際共著 0件 / うちオープンアクセス 5件）

|   |                       |
|---|-----------------------|
| 1. 著者名<br>Alastair Butler   | 4. 巻<br>13            |
| 2. 論文標題<br>From meaning representations to syntactic trees  | 5. 発行年<br>2016年       |
| 3. 雑誌名<br>Proceedings of the Thirteenth International Workshop of Logic and Engineering of Natural Language Semantics13 (LENLS13) | 6. 最初と最後の頁<br>147-160 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし  | 査読の有無<br>無            |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難  | 国際共著<br>-             |

|  |                         |
|--|-------------------------|
| 1. 著者名<br>Alastair Butler  | 4. 巻<br>10              |
| 2. 論文標題<br>DynamicPower at SemEval-2016 Task 8: Processing syntactic parse trees with a Dynamic Semantics core | 5. 発行年<br>2016年         |
| 3. 雑誌名<br>Proceedings of SemEval-2016  | 6. 最初と最後の頁<br>1148-1153 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし   | 査読の有無<br>無              |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難   | 国際共著<br>-               |

|   |                   |
|---|-------------------|
| 1. 著者名<br>Alastair Butler   | 4. 巻<br>2         |
| 2. 論文標題<br>Deterministic natural language generation from meaning representations for machine translation | 5. 発行年<br>2016年   |
| 3. 雑誌名<br>Proceedings of the 2nd Workshop on Semantics-Driven Machine Translation                         | 6. 最初と最後の頁<br>1-9 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし  | 査読の有無<br>無        |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難  | 国際共著<br>-         |

|   |                      |
|---|----------------------|
| 1. 著者名<br>周振・吉本啓  | 4. 巻<br>21           |
| 2. 論文標題<br>中国人日本語学習者のVN型二字漢語動詞の習得に関する研究: VN型二字漢語動詞の一体性の視点から | 5. 発行年<br>2015年      |
| 3. 雑誌名<br>国際文化研究  | 6. 最初と最後の頁<br>99-112 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし                              | 査読の有無<br>無           |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難                      | 国際共著<br>-            |

|  |                     |
|--|---------------------|
| 1. 著者名<br>ブラシャント・バルデシ・Alastair Butler・吉本啓・岸本秀樹 | 4. 巻<br>21          |
| 2. 論文標題<br>統語・意味解析情報付き日本語 コーパスの開発              | 5. 発行年<br>2015年     |
| 3. 雑誌名<br>言語処理学会第21回年次大会発表論文集                  | 6. 最初と最後の頁<br>20-23 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし                 | 査読の有無<br>無          |
| オープンアクセス<br>オープンアクセスとしている (また、その予定である)         | 国際共著<br>-           |

|  |                       |
|--|-----------------------|
| 1. 著者名<br>Alastair Butler and Kei Yoshimoto                      | 4. 巻<br>21            |
| 2. 論文標題<br>Large scale semantic representation with flame graphs | 5. 発行年<br>2015年       |
| 3. 雑誌名<br>言語処理学会第21回年次大会発表論文集                                    | 6. 最初と最後の頁<br>301-304 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし                                   | 査読の有無<br>無            |
| オープンアクセス<br>オープンアクセスとしている (また、その予定である)                           | 国際共著<br>-             |

|   |                       |
|---|-----------------------|
| 1. 著者名<br>Alastair Butler, Shota Hiayama and Kei Yoshimoto      | 4. 巻<br>21            |
| 2. 論文標題<br>Coindexed null elements for a Japanese parsed corpus | 5. 発行年<br>2015年       |
| 3. 雑誌名<br>言語処理学会第21回年次大会発表論文集                                   | 6. 最初と最後の頁<br>708-711 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし                                  | 査読の有無<br>無            |
| オープンアクセス<br>オープンアクセスとしている (また、その予定である)                          | 国際共著<br>-             |

|  |                       |
|--|-----------------------|
| 1. 著者名<br>周振・Alastair Butler・吉本啓   | 4. 巻<br>21            |
| 2. 論文標題<br>中国語意味解析コーパス構築のための句レベルのスコアアノテーション - 文の構成要素の間のコントロール関係の同定および否定の作用域の制御を中心に - | 5. 発行年<br>2015年       |
| 3. 雑誌名<br>言語処理学会第21回年次大会発表論文集  | 6. 最初と最後の頁<br>856-859 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし   | 査読の有無<br>無            |
| オープンアクセス<br>オープンアクセスとしている (また、その予定である)   | 国際共著<br>-             |

|  |                     |
|--|---------------------|
| 1. 著者名<br>周振・Alastair Butler・吉本啓       | 4. 巻<br>17          |
| 2. 論文標題<br>中国語結果構文の解析                  | 5. 発行年<br>2015年     |
| 3. 雑誌名<br>言語科学会第17回年次国際大会, ハンドブック      | 6. 最初と最後の頁<br>56-59 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし         | 査読の有無<br>無          |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難 | 国際共著<br>-           |

|  |                       |
|--|-----------------------|
| 1. 著者名<br>アラスティア・バトラー・吉本啓・岸本秀樹・ブラシャント・パルデン | 4. 巻<br>22            |
| 2. 論文標題<br>統語・意味解析情報付き日本語コーパスのアノテーション      | 5. 発行年<br>2016年       |
| 3. 雑誌名<br>言語処理学会第22回年次大会発表論文集              | 6. 最初と最後の頁<br>589-592 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし             | 査読の有無<br>無            |
| オープンアクセス<br>オープンアクセスとしている(また、その予定である)      | 国際共著<br>-             |

|  |                       |
|--|-----------------------|
| 1. 著者名<br>周振・Alastair Butler・吉本啓       | 4. 巻<br>22            |
| 2. 論文標題<br>中国語連体修飾節構文の解析               | 5. 発行年<br>2016年       |
| 3. 雑誌名<br>言語処理学会第22回年次大会発表論文集          | 6. 最初と最後の頁<br>809-812 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし         | 査読の有無<br>無            |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難 | 国際共著<br>-             |

〔学会発表〕 計22件(うち招待講演 0件/うち国際学会 0件)

|  |
|--|
| 1. 発表者名<br>Yusuke Kubota   |
| 2. 発表標題<br>Reconsidering the Coordinate Structure Constraint once again: Corpus-based evidence |
| 3. 学会等名<br>Conceptual and Methodological Alternatives in Theoretical Linguistics               |
| 4. 発表年<br>2018年  |

|  |
|--|
| 1. 発表者名<br>アラステア・バトラー, 長崎郁, スティーブン・ライト・ホーン, プラシャント・パルデシ, 吉本啓 |
| 2. 発表標題<br>統語解析情報付きコーパス検索用インタフェースの開発                         |
| 3. 学会等名<br>言語処理学会第24回年次大会                                    |
| 4. 発表年<br>2018年  |

|   |
|---|
| 1. 発表者名<br>Alastair Butler, Stephen Wright Horn, Iku Nagasaki   |
| 2. 発表標題<br>Seeding lexical semantics: resources using parsed corpora  |
| 3. 学会等名<br>NINJAL International Symposium "Exploiting Parsed Corpora: Application in Research, Pedagogy and Processing" |
| 4. 発表年<br>2017年   |

|   |
|---|
| 1. 発表者名<br>Hideki Kishimoto, Prashant Pardeshi  |
| 2. 発表標題<br>Parsed corpus as a source for testing generalizations in Japanese syntax.                                    |
| 3. 学会等名<br>NINJAL International Symposium "Exploiting Parsed Corpora: Application in Research, Pedagogy and Processing" |
| 4. 発表年<br>2017年   |

|   |
|---|
| 1. 発表者名<br>Yusuke Kubota, Ai Kubota   |
| 2. 発表標題<br>A case study on the Coordinate Structure Constraint in Japanese  |
| 3. 学会等名<br>NINJAL International Symposium "Exploiting Parsed Corpora: Application in Research, Pedagogy and Processing" |
| 4. 発表年<br>2017年   |

|   |
|---|
| 1. 発表者名<br>Kei Yoshimoto, Akiko Takahashi   |
| 2. 発表標題<br>Exploiting coreferential information in NPCMJ for L2 reading of Japanese texts                               |
| 3. 学会等名<br>NINJAL International Symposium "Exploiting Parsed Corpora: Application in Research, Pedagogy and Processing" |
| 4. 発表年<br>2017年   |

|   |
|---|
| 1. 発表者名<br>Alastair Butler, Stephen Wright Horn   |
| 2. 発表標題<br>Treebank Semantics parsed corpus series  |
| 3. 学会等名<br>NINJAL International Symposium "Exploiting Parsed Corpora: Application in Research, Pedagogy and Processing" |
| 4. 発表年<br>2017年   |

|                                      |
|--------------------------------------|
| 1. 発表者名<br>窪田悠介                      |
| 2. 発表標題<br>「ツリーバンク検索への「UNIX 的」アプローチ」 |
| 3. 学会等名<br>国語研究所言語資源活用ワークショップ        |
| 4. 発表年<br>2017年                      |

|   |
|---|
| 1. 発表者名<br>Alastair Butler, Ai Kubota, Shota Hiyama and Kei Yoshimoto     |
| 2. 発表標題<br>Treebank annotation of FraCaS and JSeM                         |
| 3. 学会等名<br>Logic and Engineering of Natural Language Semantics (LENLS 13) |
| 4. 発表年<br>2016年   |



|   |
|---|
| 1. 発表者名<br>Alastair Butler  |
| 2. 発表標題<br>From meaning representations to syntactic trees                |
| 3. 学会等名<br>Logic and Engineering of Natural Language Semantics (LENLS 13) |
| 4. 発表年<br>2016年   |

|   |
|---|
| 1. 発表者名<br>Alastair Butler  |
| 2. 発表標題<br>Deterministic natural language generation from meaning representations for machine translation |
| 3. 学会等名<br>2nd Workshop on Semantics-Driven Machine Translation   |
| 4. 発表年<br>2016年   |

|  |
|--|
| 1. 発表者名<br>ブラシャント・バルデシ                                 |
| 2. 発表標題<br>ワークショップ「イントロダクション」統語・意味解析情報付き日本語コーパスの構築に向けて |
| 3. 学会等名<br>日本語学会第153回大会                                |
| 4. 発表年<br>2016年  |

|  |
|--|
| 1. 発表者名<br>ブラシャント・バルデシ                                 |
| 2. 発表標題<br>ワークショップ「まとめと将来の展望」統語・意味解析情報付き日本語コーパスの構築に向けて |
| 3. 学会等名<br>日本語学会第153回大会                                |
| 4. 発表年<br>2016年  |

|  |
|--|
| 1. 発表者名<br>吉本啓   |
| 2. 発表標題<br>ワークショップ「アノテーション方式とコーパスの特色」統語・意味解析情報付き日本語コーパスの構築に向けて |
| 3. 学会等名<br>日本言語学会第153回大会                                       |
| 4. 発表年<br>2016年  |

|   |
|---|
| 1. 発表者名<br>アラスデア・バトラー、窪田愛、窪田悠介                          |
| 2. 発表標題<br>ワークショップ「デモンストレーション」統語・意味解析情報付き日本語コーパスの構築に向けて |
| 3. 学会等名<br>日本言語学会第153回大会                                |
| 4. 発表年<br>2016年   |

|  |
|--|
| 1. 発表者名<br>Alastair Butler   |
| 2. 発表標題<br>Parsed Corpus Semantics                                 |
| 3. 学会等名<br>New Landscapes in Theoretical Computational Linguistics |
| 4. 発表年<br>2016年  |

|  |
|--|
| 1. 発表者名<br>Alastair Butler, Shiro Akasegawa, Prashant Pardeshi and Kei Yoshimoto     |
| 2. 発表標題<br>A parsed corpus of Japanese enriched to reach levels of semantic analysis |
| 3. 学会等名<br>Brandeis University, Boston, USA Colloquium                               |
| 4. 発表年<br>2016年  |

|                              |
|------------------------------|
| 1. 発表者名<br>窪田悠介              |
| 2. 発表標題<br>形式意味論と計算言語学の最近の動向 |
| 3. 学会等名<br>第13回東海意味論研究会      |
| 4. 発表年<br>2016年              |

|  |
|--|
| 1. 発表者名<br>吉本啓                                   |
| 2. 発表標題<br>統語・意味解析情報を伴う日本語コーパスの開発とその日本語教育・学習への応用 |
| 3. 学会等名<br>台湾日本語言文藝研究学会第15回定例学会                  |
| 4. 発表年<br>2016年                                  |

|  |
|--|
| 1. 発表者名<br>アラスデア・パトラー、吉本 啓、岸本 秀樹、ブラシャント・バルデシ |
| 2. 発表標題<br>統語・意味解析情報付き日本語コーパスのアノテーション        |
| 3. 学会等名<br>言語処理学会 第22回年次大会                   |
| 4. 発表年<br>2016年                              |

|  |
|--|
| 1. 発表者名<br>吉本啓・ブラシャント・バルデシ                 |
| 2. 発表標題<br>文の統語・意味解析情報をタグ付けした日本語構造体コーパスの開発 |
| 3. 学会等名<br>関西言語学会ワークショップ                   |
| 4. 発表年<br>2015年                            |

|   |
|---|
| 1. 発表者名<br>Kei Yoshimoto and Alastair Butler  |
| 2. 発表標題<br>Development of Japanese Corpus Tagged with Syntactic and Semantic Information  |
| 3. 学会等名<br>The 18th Joint Workshop on Linguistics and Language Processing. Korean Society for Language and Information. Kyung Hee University, Seoul |
| 4. 発表年<br>2015年   |

〔図書〕 計0件

〔産業財産権〕

〔その他〕

-

6. 研究組織

|       | 氏名<br>(ローマ字氏名)<br>(研究者番号)                     | 所属研究機関・部局・職<br>(機関番号)   | 備考 |
|-------|---|---|----|
| 研究分担者 | 岸本 秀樹<br>(Kishimoto Hideki)<br><br>(10234220) | 神戸大学・人文学研究科・教授<br><br>(14501)                                   |    |
| 研究分担者 | 野田 尚史<br>(Noda Hisashi)<br><br>(20144545)     | 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・日本語教育研究領域・教授<br><br>(62618)           |    |
| 研究分担者 | 吉本 啓<br>(Yoshimoto Kei)<br><br>(50282017)     | 東北大学・高度教養教育・学生支援機構・教授<br><br>(11301)                            |    |
| 研究分担者 | 窪田 悠介<br>(Kubota Yusuke)<br><br>(60745149)    | 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・理論・対照研究領域・准教授<br><br>(62618)          |    |
| 研究分担者 | 長崎 郁<br>(Nagasaki Iku)<br><br>(70401445)      | 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・理論・対照研究領域・プロジェクト非常勤研究員<br><br>(62618) |    |

## 6. 研究組織（つづき）

|           | 氏名<br>(研究者番号)  | 所属研究機関・部局・職<br>(機関番号)   | 備考 |
|-----------|--|---|----|
| 研究<br>分担者 | バトラー アラスデア<br><br>(Butler Alastair)<br><br>(90588873)    | 弘前大学・人文社会科学部・准教授<br><br><br><br>(11101)                                 |    |
| 研究<br>分担者 | HORN S.W.<br><br>(Horn Stephen Wright)<br><br>(70801538) | 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・理論・対照研究領域・プロジェクト非常勤研究員<br><br><br><br>(62618) |    |
| 研究<br>分担者 | 影山 太郎<br><br>(Kageyama Taro)<br><br>(80068288)           | 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・理論・対照研究領域・その他<br><br><br><br>(62618)          |    |