

# 日本語書き言葉を対象にした省略・共参照解析の誤り分析

飯田龍<sup>1</sup> 柴田知秀<sup>2</sup> 井之上直也<sup>3</sup>

<sup>1</sup> 情報通信研究機構ユニバーサルコミュニケーション研究所

<sup>2</sup> 京都大学大学院情報学研究科 <sup>3</sup> 東北大学大学院情報科学研究科

## 1 はじめに

照応解析とは、代名詞等の指示表現が文章内で指示する表現を特定する処理をいい、指し元の表現を照応詞、指し先の表現を先行詞という。また、照応解析と類似する概念として共参照解析があり、この解析では、名詞句等の表現の対がある世界における同一の指示対象を指しているか否かを判別する処理を行う。この2つの解析の違いを例(1)と例(2)を用いて説明する。例(1)では、代名詞「それ」が前文の「プリウス」を指しており、さらに、「それ」と「プリウス」が同じ実体を指しているため、この関係を同定する処理は照応解析であり、かつ、共参照解析となる。

(1) 太郎は プリウス を買った。次の日、それ に乗って会社へ行った。

一方、例(2)では「その車」が「プリウス」を指しているため、この関係を同定する処理は照応解析となるが、1文目の「プリウス」は「太郎が買ったプリウス」であるのに対し、2文目の「その車」は「次郎が買ったプリウス」であるため、この2つの表現の間には共参照関係が成り立たない。このため、この2つの文の中に共参照解析の問題として解くべき関係は存在しないことになる。

(2) 太郎は プリウス を買った。次郎も その車 を買った。

また、照応解析の問題のうち、特に照応詞が省略された場合が存在する。この省略された照応詞をゼロ代名詞と呼び、このゼロ代名詞に限定して照応解析を行う問題をゼロ照応解析、もしくは省略解析と呼ぶ。例えば、例(3)の2文目の「乗って」や「行った」のガ格は省略されており、前文に出現している「太郎」を先行詞として同定する必要がある。

(3) 太郎 はプリウスを買った。次の日、(φガ) それに乗って (φガ) 会社へ行った。

また、例(3)の文を越えて先行詞を特定する問題に加え、同一文内に複数の述語が出現し、その述語の項が述語と直接の係り受け関係にない場合には、述語の格要素が省略されているとみなし、省略解析の問題として扱う場合がある。例えば、例(4)では「太郎がプリウスを買う」、

「太郎がそれに乗る」、「太郎が会社へ行く」という述語項関係が成り立つが、この文では「太郎」と「買い」が係り受け関係にある場合、それ以外の述語「乗る」「行く」のガ格は省略されているとみなし、文内の省略解析として「乗る」「行く」の省略された格要素の先行詞が「太郎」であることを同定する必要がある。

(4) 太郎 はプリウスを買い、次の日、(φガ) それに乗って (φガ) 会社へ行った。

さらに、照応解析では先行詞が必ずしも文章中に出現するとは限らない。例えば、例(5)では、ゼロ代名詞「φ」がこの文の書き手であり、「その花」が書き手の世界におけるある花であるが、その指し先が明示的にテキスト中に出現していない。

(5) その花 は枯れそうだったので (φガ) 水をやった。このような場合は、ゼロ代名詞「φ」の指し先がテキスト中に存在しない一人称であるという外界照応の関係を特定する問題として扱われている [19, 21, 15]。

日本語を対象にした照応・省略解析は、その研究の初期段階では人手で作成された規則に基づいて関係を同定する手法 [7, 8, 20] が提案され、その後は他の研究分野と同様に、コーパスアノテーションと機械学習に基づく解析手法が提案されている [1, 2, 9, 5, 18]。このコーパスアノテーションに基づく解析手法では、京都大学テキストコーパス 4.0 [12] や NAIST テキストコーパス [24]、プログコーパス [13]、Web コーパス [22] 等、さまざまなコーパスが利用されており、さらに、今後日本語の基盤技術開発のための共通データとなることが予想される日本語書き言葉均衡コーパス [11] に対しても照応・省略関係のアノテーションが進められている。これらのコーパスで採用されている省略関係のアノテーションの仕様は、前述の例(4)や例(5)のような文内や文間の省略関係をアノテーションするという共通点はあるが、例えば、京都大学テキストコーパスが受け身や使役等まで含めた述語の表層格にアノテーションしているのに対し、NAIST テキストコーパスで述語の原形の格にアノテーションしているといった細かな違いが存在する。つまり、異なる仕様で作られたコーパスに対して、いくつかの解析器が

開発されているという状況にあるが、現状ではこの違いが影響するほど解析器の性能が高くないと考えられるため、本研究で行う誤り分析では、解析器の学習時の述語認定の基準と、分析対象とする述語認定の基準が若干異なっている、その違いは無視して解析器の出力結果を分析することとする。我々の行った誤り分析では、省略・共参照解析ツールとして日本語構文解析器 KNP<sup>1</sup> と述語項構造解析器 SynCha<sup>2</sup> を用いた。

本稿では、まず、本分析で利用する 2 節で KNP と SynCha の概要を説明し、次に、3 節で分析対象とするコーパスについて説明する。4 節と 5 節で省略解析と共参照解析の分析結果をそれぞれ報告し、最後に 6 節でまとめと今後の課題について議論する。

## 2 分析に利用するツール

前述のように、照応・共参照の誤り分析のために、日本語構文解析器 KNP と述語項構造解析器 SynCha を利用する。KNP の省略解析の実体は笹野ら [18] の識別モデルであり、大規模な格フレームを利用して格要素を特定しているという特徴がある。また、KNP の共参照解析は、基本的には文献 [17] で導入された共参照関係認定基準 1「照応詞、先行詞候補を含む文節を比較し、照応詞候補を含む文節の照応詞候補以前の部分が、先行詞候補を含む文節に含まれている場合、同一、または、簡潔であるとする」にしたがって共参照関係を同定している [16]。

一方、SynCha の省略解析では、文献 [4] のガ格の省略解析手法をガ格以外のヲ格、ニ格にも適用するよう修正された解析手法が実装されている。この手法では、省略検出モデル、文内先行詞同定モデル、文間先行詞同定モデル等の出力結果を整数計画法に基づいて最適化することで最終的な省略解析の結果を得ている。また、SynCha の共参照解析では、基本的には文献 [4] の解析手法を採用しているが、単純に照応詞候補と先行詞候補との組み合わせを解析対象とすると、解析に時間がかかるため、文献 [3] のキャッシュモデルに基づいて、文単位で先行詞候補の選別を行っている。

## 3 分析対象のデータ

誤り分析に利用するコーパスとして、プロジェクト「ロボットは東大に入れるか」[6] で分析対象の一つとされている世界史 B の教科書内のテキストの一部と、日本語書き言葉均衡コーパス (BCCWJ) [11] の一部にアノテーションされた述語項関係・共参照関係の情報を利用する。以降、前者を教科書コーパス、後者を BCCWJ

コーパスと呼ぶ。教科書コーパスでは一貫した表記で記載され、教科書という媒体の性質上、読み手に理解しやすいような形式で文章が記述されているため、共参照・省略解析の対象としては比較的容易に解析できるコーパスであると考えられる。一方、BCCWJ コーパスには書籍や白書等多様な文章スタイルが含まれ、文章によっては照応・省略関係の同定が人間でも困難な事例が含まれるため、教科書コーパスと比較して解析が困難になることが予測される。

教科書コーパスを利用する際は、教科書の本文を大見出し単位で分割し、そこからランダムに 41 個抽出したものに対して、述語項関係・共参照関係をアノテーションしたものを利用する。アノテーションの基準としては概ね NAIST テキストコーパスのアノテーション仕様に準拠してアノテーションを行ったが、このアノテーションでは述語の原形ではなく、格交替をとまなう場合は交替まで含めた対象、例えば、「考えられる」や「導入させる」といった表現に対してガ格、ヲ格、ニ格をアノテーションしている点や、外界照応についてもアノテーションされていない点が NAIST テキストコーパス内のアノテーション結果とは異なっている。一方、BCCWJ に対するアノテーションでは、述語の原形に対する述語項関係をアノテーションする等、NAIST テキストコーパスのアノテーション仕様に準拠してアノテーションが行われている。BCCWJ にはさまざまなレジスタの情報が含まれているが、本稿で対象とする誤り分析のために、BCCWJ コアデータ中の書籍レジスタから 6 テキスト、新聞レジスタから 4 テキスト、知恵袋レジスタから 6 テキスト、白書レジスタから 2 テキストをランダムに抽出して利用した。各コーパスの述語項構造関係のアノテーション個数を表 1 にまとめる。省略関係のアノテーションの個数を求めるために文の係り受け構造が必要となるが、この構造を得るために、係り受け解析器 CaboCha<sup>3</sup> が出力する係り受け構造を利用した。BCCWJ コーパスでは外界照応がアノテーションされているが、本分析で利用する 2 つのツールはともに外界照応の関係を同定できないため、外界照応の関係は分析の対象外とする。

また、各コーパスにアノテーションされた共参照関係の個数としては、教科書コーパス中に 290、BCCWJ コーパス中には 168 の共参照の連鎖<sup>4</sup>が存在し、それぞれの連鎖に含まれる共参照関係となる表現の個数の平均は教科書コーパスでは 3.45、BCCWJ コーパスでは 2.96 であった。

<sup>1</sup>4.12 版を利用。

<http://nlp.ist.i.kyoto-u.ac.jp/index.php?cmd=read&page=KNP>

<sup>2</sup>0.3.1 版を利用。 <https://sites.google.com/site/ryuuiida/syncha/>

<sup>3</sup><https://code.google.com/p/cabochoa/>

<sup>4</sup>同一実体を指す表現の集合をここでは 1 つの共参照連鎖とみなしている。

表 1: アノテーションされた述語項関係の個数

		係り受け	省略	省略 (文内)	省略 (文間)	省略 (外界)	合計
教科書	ガ格	3,815	2,221	2,009	212	-	6,036
	ヲ格	2,378	241	240	1	-	2,619
	ニ格	1,454	165	160	5	-	1,619
BCCWJ	ガ格	1,782	2,054	831	581	642	3,836
	ヲ格	1,235	330	232	71	27	1,565
	ニ格	504	192	123	47	22	696

表 3: KNP と SynCha が正解の事例 (教科書コーパス)

	省略 (文内)	省略 (文間)	省略 (文内+文間)	合計
ガ格	309	7	302	618
ヲ格	0	0	0	0
ニ格	1	0	0	1
合計	310	7	302	619

表 4: KNP と SynCha が正解の事例 (BCCWJ コーパス)

	省略 (文内)	省略 (文間)	省略 (文内+文間)	合計
ガ格	66	6	56	128
ヲ格	0	0	0	0
ニ格	0	0	0	0
合計	66	6	56	128

## 4 省略解析の分析

KNP と SynCha を用い、3 節で導入した 2 つのコーパス中の省略関係をどの程度解析できるかを調査するために、再現率、精度、F 値で評価した結果を表 2 にまとめる。この結果より、どちらのツールも教科書を対象にした場合は F 値が約 4 割を上回る結果を得ているのに対し、BCCWJ を対象にした場合は F 値が 2 割未満となる結果を得ていることがわかる。この結果からわかるように、教科書コーパスのほうが簡単に解ける問題が多く含まれていることがわかる。以降の分析では、教科書コーパスと BCCWJ コーパスのそれぞれに対して、KNP と SynCha がともに解析できた事例、また、どちらのツールでも解析できなかった事例の一部を手で分析し、どのような特徴があるのかを明らかにする。

### 4.1 教科書コーパスで省略関係を同定できた事例

まず、教科書コーパスを対象に 2 つのツールがともに解析できた事例を調査した結果を報告する。正解できた事例集合のうち、先行詞がどのような出現位置に出現したか (ゼロ代名詞と同一の文内に出現 (文内)、ゼロ代名詞と異なる文内に出現 (文間)、同一文内にも異なる文にも先行詞が出現 (文内+文間)) を格ごとに調査した結果を表 3 にまとめる。この結果から、解析できた事例のほとんどが文内のガ格であることがわかる。解析できた 619 事例からランダムに 100 事例を選択し、手で事例を調査し、どのような特徴を持つかを調べた。この結果、100 事例のうち 91 事例は、例 (6) のように助詞「は」で主題化された名詞句が同一文内でその名詞句の後に出現する述語のガ格のゼロ代名詞の先行詞となる場合であった。

(6) 原人<sub>1</sub> は、アフリカを ( $\phi_1$  ガ) 出て、各地に分散した。これ以外には、例 (7) のように助詞「は」ではなく助詞「が」をともなう名詞が先行詞となる場合が 3 事例、例 (8) のような述語より後に先行詞が出現する場合は 3 事例、例 (9) のように文間照応となる場合は 2 事例、係り受け関係の誤りが 1 事例であった。

(7) すでに半島に南下していたインド=ヨーロッパ語系の定住民のなかから、やがてラテン系のローマ人<sub>1</sub> が勢力をのぼしてエトルリア人の勢力を退け、山岳地域の諸部族を ( $\phi_1$  ガ) やぶり、前 3 世紀初頭には半島全域に覇権を確立した。

(8) 「海の民」の襲撃や新たな民族移動による激動期を ( $\phi_1$  ガ) 迎え、以後 400 年ほどの間、鉄器時代に移行しながら ギリシア本土<sub>1</sub> は停滞と混乱の時代となった。

(9) ギリシア人<sub>1</sub> はポリスごとに対立・抗争をくりかえし、一つの勢力にまとまることはなかった。しかし、共通の言語を ( $\phi_1$  ガ) もち、オリンピアの祭典やデルフォイの神託を通じて、同一の民族としての自覚を失うことはなかった。

### 4.2 BCCWJ コーパスで省略解析を同定できた事例

次に、BCCWJ コーパスについても 2 つのツールがともに正しい省略関係を同定した事例について調査した。結果を表 4 にまとめる。この結果より、教科書コーパスと同様に、ツール共通に解析できる事例はガ格のみであることがわかる。解析できた事例のうち、文間の省略関係も同定できた例が数例存在するが、これらは例 (10) のように、同一文内で「服は」のように他の主題化された表現は存在するものの、格フレームもしくは述語と項の相互情報量といった情報で比較的容易に項とならないことがわかり、かつ、近傍に正解となる項も「は」で主題化されている場合のように、解析が比較的容易な場合に限定されている。

(10) 私<sub>1</sub> は毎日施設に通って、シスターやボランティアの女の子たちと一緒に洗濯をした。服は毎日 ( $\phi_1$ ) 洗っていたから、百人分の量は大変なものだ。

さらに、正しく解析できた事例のうちランダムに抽出した 100 事例を手で分析する。この分析では主に先行詞に後続する助詞に関して分類を行った。結果を表 5 に示す。この表からもわかるように、先行詞として選ばれ

表 2: SynCha と KNP の解析結果

	SynCha								KNP							
	教科書				BCCWJ				教科書				BCCWJ			
	ガ	ヲ	ニ	合計	ガ	ヲ	ニ	合計	ガ	ヲ	ニ	合計	ガ	ヲ	ニ	合計
R	0.448	0.022	0.011	0.401	0.175	0.016	0.00	0.139	0.511	0.200	0.237	0.479	0.217	0.075	0.119	0.190
P	0.422	0.032	1.000	0.406	0.236	0.070	0.00	0.225	0.455	0.075	0.083	0.368	0.219	0.063	0.083	0.180
F	0.435	0.026	0.017	0.404	0.201	0.026	-	0.172	0.481	0.109	0.123	0.416	0.218	0.069	0.097	0.185

表 5: KNP と SynCha が正解した事例の特徴

出現位置	特徴	事例数
文内照応	先行詞が助詞「は」をともなう	73
	先行詞が助詞「が」をともなう	10
	先行詞が助詞「も」をともなう	1
	先行詞が助詞「では」をともなう	3
	先行詞が助詞「には」をともなう	1
文内照応 (述語の後)	先行詞が助詞「は」をともなう	4
	先行詞が助詞「が」をともなう	2
文間照応	助詞「は」をともなう	3
	助詞「が」をともなう	1
その他	係り受け解析誤り	2

表 6: 省略解析で KNP と SynCha が先行詞を出力するが不正解の事例 (教科書コーパス)

	省略(文内)	省略(文間)	省略(文内+文間)	合計
ガ格	84	28	26	138
ヲ格	10	1	1	12
ニ格	5	0	1	6
合計	99	29	28	156

た名詞句のほとんどは「は」で主題化されていることがわかり、それ以外の解析が困難な場合 (例えば「A の B」の A を先行詞として選択する場合等) は含まれていないことがわかった。

#### 4.3 教科書コーパスで省略関係を同定できなかった事例

次に、教科書コーパスを対象に、KNP と SynCha がともに解析を誤った事例を調査する。誤り方の種類はツールが先行詞を出力したが、その出力が不正解である場合と、正解となる省略関係を出力できない誤りの 2 つに分類できる。各ツールの出力がともに前者に該当した場合の個数を表 6 に、後者に該当した場合の個数を表 7 にまとめる。

この 2 つの表から、どちらの場合もガ格以外の格についての誤りが多く含まれていることがわかる。また、2 つの表を比較してみると、解析結果を出力して誤る場合よりも、出力できずに誤る場合が多いことがわかる。これは省略解析の評価尺度として F 値が利用されていることが起因していると考えられる。F 値に基づく評価では、アノテーションされた正解の省略関係に対して、誤った関係を出力した場合は再現率と精度をともに低下させるのに対し、結果を出力しない場合は再現率しか低下させないため、解析が困難な事例については積極的に関係は出力しないほうが結果的に F 値を低下させないため、このような結果を得たと考えられる。F 値に基づいて省略解析が評価される場合は、信頼できる結果だけを出力し、それ以外の解析が誤る可能性のある関係は出

表 7: 省略解析で KNP と SynCha が先行詞を出力できず不正解の事例 (教科書コーパス)

	省略(文内)	省略(文間)	省略(文内+文間)	合計
ガ格	178	46	56	280
ヲ格	56	0	8	64
ニ格	54	5	9	68
合計	288	51	73	412

表 8: 省略解析: 先行詞を出力したが不正解になった事例の人手分類 (教科書コーパス)

特徴	事例数
正解が文内、ツールの出力は文間	1
正解が文間、ツールの出力は文内	14
正解が連体修飾の亜種、システムが文間	1
その他	11
アノテーションの誤り・問題	17
正解が助詞「が」、ツールの出力が助詞「は」	4
正解は助詞「は」、ツールの出力が助詞「が」	1
丸括弧が影響	20
正解は助詞「は」、ツールの出力が助詞「は」	13
述語が「A の B」に係る	11
述語が「AB」に係る	1
その他	5
合計	55

力しないという戦略が結果的に良い F 値を得ることになるため、この偏りを解消するためには、例えば、F 値以外の評価尺度として、コーパス中の省略関係を特定するための難度を考慮した評価尺度等の導入を考える必要がある。

#### 4.3.1 ツールが先行詞を出力したが不正解となった事例の分析

次に、2 つのツールが省略関係を出力して誤った 156 事例からランダムに 100 事例を抽出し、先行詞の出現位置に応じて人手で誤りを調査した結果を表 8 示す。この際、文内の省略関係に対してツールも文内の先行詞候補を選択して誤った事例が頻出したため、この場合についてはさらに細かく分類した。

表 8 に示した誤りの分類のうち、主要な誤りを以下にまとめる。まず、14 件存在した正解が文間照応であったにも関わらず、文内照応の関係を出力した事例について説明する。これは、例えば、例 (11) のように KNP と SynCha で出力する傾向が異なるため、誤りの原因の一般的な特徴付けは困難だが、各ツールの誤りの傾向を見る限り、典型的には KNP が連体修飾の関係として誤った項を選択するが故に、省略解析の問題を解く段階に移行できないという問題があるのに対し、SynCha では過剰に近傍の「は」で主題化された名詞句を先行詞として選択してしまうという傾向があることがわかった。

(11) オスマン帝国<sub>ans</sub> は、こうした軍事封土制(ティマール制)によって軍事力を確保したほか、バルカン半島のキリスト教徒の子弟を徴用し、イスラーム教の教育をしたうえで、スルタン直属の官僚としたり、常備軍(イエニチェリ)に編成した。いっぽう、エジプト<sub>syn</sub> やイラクなど新しく(φガ)獲得した領土<sub>knp</sub> は、総督や現地の支配層に統治をゆだね、貢納金の支払いを義務づけた。

このような各ツールの解析の傾向に加え、省略解析の手法に共通する特性として、文間照応の候補よりも文内照応の候補を出力しやすい傾向にあることがわかる。これは、文間の省略関係よりも文内の省略関係の同定性能が高いため、文間の候補も先行詞同定の対象に含まれている場合であっても、文内のそれらしい先行詞候補を先行詞として出力してしまうためだと考えられる。

上述の推測を含め、ツールが出力する結果の分析は、解析時に利用される特徴のうち、どれに着目して考察するかで分析結果が変わってしまうため、ツールの出力する結果の分析は困難を極める。ただし、ツールが出力する結果は省略解析の問題の性質上、述語の選択制限を満たすか否かという観点からは比較的容易に分析できると考えられる。例えば、例(12)では、「滅ぼす」のガ格として「時代」が選ばれているが、これはこの文脈においては選択制限を満たしていない。

(12) 西ローマ帝国の滅亡にもかかわらず、東ローマ帝国<sub>ans</sub> の活況はかわらなかった。首都コンスタンティノープルをはじめとするヘレニズム時代にさかのぼる諸都市は繁栄をつづけ、商工業もさかんであった。6世紀のユスティニアヌス帝の時代<sub>knp,syn</sub> には、東ゴート王国とヴァンダル王国を(φガ)滅ぼしてイタリア半島とアフリカ北岸を奪回し、西ゴート王国からイベリア半島南部を奪って、帝国領の大半を回復した。

このように、各ツールが出力した結果のうち、選択制限を満たした出力が得られる割合を調べることはそれほど難しくないと考えられるが、ツールが出力する結果は複合的にさまざまな特徴を組み合わせることで出力しているため、選択制限そのものを評価することにはつながらない。そこで、ここでは選択制限の評価問題について提案したい。例えば、例(12)では「滅ぼす」の選択制限を満たす名詞句として「西ローマ帝国」や正解となる「東ローマ帝国」等、複数の選択選好を満たす名詞句と、「時代」のような選択制限を満たさないと考えられる名詞句が存在するが、このうち、選択制限を満たす名詞句全てをアノテーションし、そのアノテーション結果を用いて、省略解析器で利用されている格フレーム辞書や、述語と項の

共起に基づく選択選好を評価することで、省略解析器で用いるべき選択制限・選択選好の言語資源の評価・選別を行えると考えられる。このような省略解析の個別の特徴の調査については、今後十分な規模の評価用データを作成し、その上で評価を行いたい。

また、表8に含まれる誤り事例には、先行詞の候補がガ格として補完できるスコープを越えた箇所から選択される例も存在した。例えば、例(13)では、「フランス」はその出現位置よりも後で主語となることが予想されるのに対し、その表現よりも前のゼロ代名詞の先行詞として補完されることでおかしな結果となっていることができる。

(13) [ところが、神のお告げを受けたと信じる農民の娘ジャンヌ＝ダルク<sub>ans</sub> が先頭に立ってオルレアンを攻囲を(φガ)やぶると、IIフランス<sub>knp,syn</sub> は反攻に転じ、シャルル7世のもとで、カレー市を除く全領土を確保して、戦争は終わった。]<sup>5</sup>

このようなガ格のスコープについても、上述の選択制限の分析と同様に、全自動もしくは人手でスコープをアノテーションを行い、その結果を利用してスコープを越えた先行詞を同定した場合になぜそれが起こるのかを分析することで新たな知見が得られると考えられる。

さらに、後述する名詞句共参照関係の分析でも見られるように、それ自体では具体的に何を指しているかの判断が困難な表現を明示的に扱う必要がある。例えば、例(14)の「中心」はその名詞単体では何の中心であるかわからないため、述語の選択制限を満たしているかの判断ができない。この「中心」は厳密には、「交易活動の中心」と解釈でき、その結果を利用して「交易活動」が「(交易網を)握る」という選択制限を満たしていることの判断が可能になる。

(14) シリアを故地とするフェニキア人<sub>ans</sub> は西地中海の各地に植民活動を行い、とくに北アフリカを拠点としたカルタゴはフェニキア系植民市群の頂点に立った。その交易活動の中心<sub>knp,syn</sub> は、金、銀、錫、銅を独占することになり、前6世紀ごろから西地中海の交易網を(φガ)握って大勢力を誇った。

この「中心」や「ほとんど」のような名詞句は有限であると考えられるため、その表現のリストを作成し、それを共有することが重要だと考えられる。

#### 4.3.2 ツールが先行詞を出力できず不正解となった事例の分析

次に、2つのツールが省略関係を出力できなかった場合について、100事例を人手で調査し、誤り事例を分類した結果を表9にまとめる。全般的な傾向として、この

<sup>5</sup>[ ]は説明のために著者が追記した。

表 9: 省略解析: 先行詞が出力できず不正解になった事例の人手分類 (教科書コーパス)

特徴	事例数
アノテーションの誤り・問題	15
機能語相当表現へのアノテーション	10
「いわれた」のような外界照応の問題と混在	9
名詞+”だ”の格要素	9
離れた位置 (1 文前 or 2 文前) に先行詞が出現	6
ガ格で先行詞が述語より後に出現	5
名詞句チャンキングの誤り	4
丸括弧の問題	4
ニ格で先行詞が述語より後に出現	3
ひらがな表記が影響	2
ニ格の解析誤り	2
文末の名詞句が先行詞となる	2
その他	29

表にまとめられた事例の多くはアノテーションの揺れやその結果起こる問題、もしくはアノテーションの仕様のためにアノテーションされにくい現象に該当することがわかった。例えば、表 9 のうち最頻出の問題はアノテーションの誤りであるが、この一部はアノテーションの仕様のためにアノテーションされた関係を捉えられないものである。例えば、(15) では、初出の「ギリシア人」が正解としてアノテーションされており、二つ目に出現した「ギリシア人」がツールの出力である。この際、2 つの「ギリシア人」が共参照関係としてアノテーションされている場合は、その情報参照して 2 つ目の「ギリシア人」も正解とすることができるのだが、現在採用している共参照認定の基準では 2 つの名詞句がある世界において厳密に一致していることを共参照関係を特定する条件としているため、判断が悩ましい場合や総称的な表現の間の共参照関係はアノテーションされない傾向にあり、結果的に例 (15) のような場合は不正解となってしまう。

#### (15) ギリシアの古典文明

ギリシア人<sub>ans</sub> はオリエントの先進文明を受け入れつつ、人間中心の考え方にもとづいて合理的な精神をつちかひながら、独自の文明を生みだした。

ギリシア人<sub>syn</sub> の心にはオリンポス 12 神を中心とする神話の世界が生きていた。人間の姿をした神々はそれぞれが豊かな個性をもち、喜怒哀楽をあらわに人間にはたらきかけると (φ ガ) 考えられた。

この問題を回避するために、省略解析のアノテーション側で網羅的に項となる表現をアノテーションすることも考えられるが、テキストが長くなればなるほど、その負荷は大きくなるので、萩行ら [22] が行っているような、テキストの短い抜粋をアノテーションの対象とし、その中で網羅的に省略関係をアノテーションすることも考えられる。ただ、このような特殊な場合を除いてもアノテーション誤りの全体に占める割合は大きい。十分な規模のデータを作成するには、典型的には外注に頼らざるを得ないため、この結果、研究者側は作業内容の確認や、

その内容の修正を頻繁に行うことが難しくなってしまう。その結果、作業結果に少なくないアノテーションの誤りが含まれることを回避するのは難しい。十分なトライアルを行ったとしても、省略・共参照関係のアノテーションでは頻繁に想定外の問題が発生するため、その想定外の場合に行う作業を作業者の判断に任せることになってしまう。このような問題を回避するには、コストをかけてでも 1 つのテキストに複数人でアノテーションを行い、その結果揺れが発生した箇所を細かく修正する、もしくは、光田ら [14] が研究を進めているような個々のアノテーション事例の誤りの可能性を推定するような技術を同時に研究開発する必要があると考えられる。

次に頻出した問題が「～を通じて」や「～に応じて」のような機能的な表現とみなせる句に対してガ格をアノテーションしているために、それをツール側が出力しないという問題である。このような機能的な表現に対するアノテーションは作業者に周知するものの、どのような場合に付けるかは、例えば、「加える」が実際にモノを加える意味で出現している場合にはアノテーションするが「～に加え」のような句で出現している場合にはそれを除外するという判断が必要になり、そのような曖昧性が作業者の負荷を高くしているのだと考えられる。「という」、「といわれる」、「と考えられる」といったテキスト内の表現をガ格とすべきか、それとも「我々」のような主体をガ格に埋めるかの判断が揺れ、結果的にアノテーションに揺れが生じるという問題も起こっていることも確認できた。同様に「～するものだ」の「ものだ」にアノテーションされているが、それをツール側が出力しないという問題もこれに類する。

これらアノテーションに関係する問題に加え、ツールが述語独立に問題を解くことで起こる問題も存在する。例えば、例 (16) では「生長する」、「実をつける」、「枯れる」のガ格はすべて「(冬型の) 植物」である。この例で「枯れる」は係り受け関係にある「(冬型の) 植物」との連体修飾関係を判別するが、それ以外の動詞はその連体修飾関係とは独立に項を探索するために、問題が難しくなる。

#### (16) そこでは、冬雨で (φ<sub>1</sub> ガ) 生長し 春に実をつけ、夏の乾燥で枯れる冬型の 植物<sub>1</sub> が生育する。

このような問題を回避するために、項の共有関係を事前に捉え、その結果を利用して共通の項を捉える必要があると考えられる。

また、これ以外にも文末に出現する名詞句を先行詞とする場合や、名詞句のチャンキングを誤ることで先行詞となる表現を捉えられない場合、「かくす」や「こうむる」といった普段ひらがな表記で書かれない動詞がひら

表 10: KNP と SynCha が先行詞を出力するが不正解の事例 (BCCWJ コーパス)

	省略 (文内)	省略 (文間)	省略 (文内+文間)	合計
ガ格	54	98	13	165
ヲ格	5	1	0	6
ニ格	3	1	1	5
合計	62	100	14	176

表 11: KNP と SynCha が先行詞を出力できず不正解の事例 (BCCWJ コーパス)

	省略 (文内)	省略 (文間)	省略 (文内+文間)	合計
ガ格	154	274	40	468
ヲ格	102	63	18	183
ニ格	59	45	19	123
合計	315	382	77	774

がなで書かれることで出現表記に基づいて言語資源を作成している場合にそれが利用できなくなるといった個別の問題も存在した。

#### 4.4 BCCWJ コーパスで省略関係を同定できなかった事例

BCCWJ コーパスに対しても、前節と同様に、先行詞を出力したが解析を誤った場合と先行詞を出力できなかった場合についてその個数を調査した結果をそれぞれ表 10 と表 11 にまとめる。この 2 つの表より、教科書コーパスと同様に、省略関係を出力せずに誤っている個数が多いことがわかる。この傾向は特にガ格以外の格で顕著で、ヲ格とニ格の省略のほとんどは出力されていないことがわかる。また、誤りの傾向としては概ね教科書コーパスと同じ傾向が見られたため、詳細は割愛する。

### 5 共参照解析の分析

3 節に示した 2 つのコーパス中の共参照連鎖に対し、予備調査として SynCha を適用した結果の調査を行った。この調査では、SynCha が出力した共参照連鎖が具体的にどのようなものであるかを調査する。まず、自動解析の結果得られた共参照連鎖に含まれる表現の集合を、その出現順で対にし、その対に含まれる表現がどのようなものであるかを調査した。この結果、ほとんどの対が完全に文字列一致している表現対となっていることがわかった (教科書コーパスで 98%、BCCWJ コーパスで 99%)。残りの対についても、部分的に文字列が一致しており、これらの文字列一致で自動的に出力された共参照関係の多くは人名や組織名、地名等の固有表現であることがわかった。これらの表現については偶然同名で同じ記事内にその表現が出力されない限り、ほぼ間違いなく共参照関係となるため、新聞記事コーパスを利用して学習して構築された共参照解析モデルはこのような表現対を特に共参照関係として出力しやすいことがわかる。また、この文字列一致情報に過剰に頼る解析の副作用として「それ」と「それ」のような代名詞を共参照関係として出力したり、「彼」と「彼女」のような確実に同じ実際を指さない表現対も共参照関係として同定してしまっている

表 12: 教科書コーパス中の共参照関係の分類結果

特徴	個数
一方が代名詞	12
括弧内の表現が影響	11
文字列の部分一致	7
別名	7
一方が指示連体詞+名詞	5
データ作成・タグ付けの誤り	4
特殊な指示表現	2
ほぼ等価な表現	1
特殊な先行詞	1

ため、普通名詞についてはその表現の出現文脈や名詞句の修飾表現に基づいた共参照関係とならない場合の判別が必要となる。

このように SynCha は新聞記事から共参照関係を学習しているものの、アノテーションされた関係のほとんどが文字列一致の特徴を持つため、出力した結果のほとんどは完全に文字列が一致するという結果になっており、また、KNP についても 2 節で述べたように、文字列一致を前提として解析するため、文字列が一致しない場合を考慮できないことになる。

そこで、本分析では、2 つのコーパス中でアノテーションされた共参照関係のうち、現状では解析できない文字列が一致しない場合、特に完全一致しない共参照関係の中からそれぞれランダムで 50 事例ずつを抽出し、それらの特徴について調査した。

#### 5.1 教科書コーパスの共参照関係の分析

教科書コーパスに含まれる共参照関係となる表現の対のうち、文字列が完全一致しない 555 事例から 50 事例をランダムに抽出し、その事例を特徴ごとに分類した結果を表 12 にまとめる。以降で、それぞれについて説明する。

表 12 から、文字列が一致しない共参照関係のうち最も頻度が多い特徴として一方が代名詞となっていることがわかる。これは教科書テキストの特徴として、歴史上の人物や事件の説明中に、その人名や事件を指示表現で指すことが多いために、この関係が多く含まれることになる。例えば、例 (17) の「改革運動」が「それ」で指される、例 (18) の「カエサル」が「彼」で指される関係がこれに対応する。

(17) 1 世紀になると、クリュニー修道院を先陣として 改革運動 がおこり、それを背景として教皇グレゴリウス 7 世は厳格な規律を求める大改革に着手した。

(18) しかし、三者の協力関係は長くはつづかず、ガリア遠征で勢威を高めた カエサル は、ポンペイウスを倒し、前 46 年には事実上の独裁者となった。彼 は、政治や社会の安定のために諸改革を断行したが、前 44 年、共和政信奉派の手で暗殺された。

このような代名詞の利用は、これまで主に研究対象とさ

れてきた新聞記事コーパスでは、その記事と短く簡潔に記述するという特徴から省略されることが多く、コーパス内にほとんど出現しないため、それらを分析・学習した解析器では解析できない。今後は BCCWJ 等の新聞記事以外のコーパスが研究対象として広く利用されることで、この代名詞が問題に含まれないという問題は解消されることが考えられる。

次に、代名詞の利用に次いで、括弧の利用が共参照関係の同定に影響していることがわかった。これは、例(19)の「ジズヤ」が「(ヒンドゥー教徒に課せられていた)人頭税」と同じであるために共参照関係としてアノテーションされたことに影響している。ある名称とその説明の表現を共参照関係とアノテーションするかどうかについては議論の余地はあるが、例(20)のように、別名として記述する場合には共参照関係とすべきだと考えられる。

(19) 彼はまた、ヒンドゥー教徒のラージプート諸侯と婚姻関係を結んで、彼らを軍に組みこみ、また、ヒンドゥー教徒に課せられていた人頭税(ジズヤ)を廃止するなど、ヒンドゥー教徒に対して融和政策をとった。

(20) テイ和は チャンパー (占城) とマラッカ海峡に面した港市国家マラッカを根拠地にして、東南アジア、インド洋各国に中国への朝貢をすすめた。

括弧の扱いについて省略関係も関係するため、コンセンサスが得られるような問題の定義についての議論が必要となる。

また、アノテーションされた共参照関係の中には同一文章内で「イギリス」と「英」、「ジラフ」と「キリン」のような別名で同じ実体や実体のクラスを指すことがあり、これらも共参照関係の同定対象となる。このような別名を対象にする場合には、基本的にはその別名に関する知識がなければ解析できないと考えられるが、「ジラフをキリンというのはこれにはじまる」や「ハンは漢字で可汗と書き、カン(カガン)と読むこともある」のように対象文章中にそれを捉える手がかりが記述されることもあるため、そのような説明が出現したときにそこから頑健に別名を捉えるという方向性も考えられる。

さらに、共参照関係の厳密性についても議論の余地がある。例(21)では、「ヨーロッパ諸国」と「ヨーロッパ」はその2つの表現の出現文脈からほとんど同じ実体を指しているとみなして良いため、共参照関係としてアノテーションされているが、このような関係をもし認めるのであれば、どこまでを許容するのかを決める必要がある。新聞記事においても「日本」と「日本政府」がほぼ等価な意味で使われることが頻繁に存在するため、新聞

表 13: BCCWJ コーパス中の共参照関係の分類結果

特徴	個数
一方もしくは両方が代名詞	24
言い換え、(人物の)役割、別名の認識が必要	9
タグ付け・データ作成の誤り	6
会話文の中に出現することが解析に影響	5
事態(節)を指す	3
一方が指示連体詞+名詞	2
ほぼ等価な表現	1

記事コーパスに共参照関係を一貫してアノテーションする場合でも問題になる。

(21) ヨーロッパ諸国の貿易が活発になるなか、17世紀前半には独立後の勢いをもったオランダが、...武力をも用いて権益の拡大をはかると、おりからの胡椒価格の暴落もあって、オランダの商業覇権は急速に傾いていった。

17世紀から18世紀にかけ、インド綿布はヨーロッパでの評価が高まり、奴隷貿易の代価にも用いられる国際商品になった。

最後に、代名詞については「それ」や「彼」のような典型的な表現に加え、それ以外にも「前者」、「後者」、「もう一方」、「残りの2つ」等が存在し、それらを適切に捉えて指し先を特定する必要がある。同様に、先行詞についても先行詞となる最右の形態素が「ほとんど」「一部」のように、それ自体ではほとんど情報を持たない表現があり、それを先行詞の候補とする場合には、「(騎馬遊牧民の)ほとんど」が「騎馬遊牧民」の意味的な特性を継承して解析時に利用される必要がある。このような表現の網羅と適切な取り扱いも今後の課題となる。

## 5.2 BCCWJ コーパスの共参照関係の分析

教科書コーパスを分析した場合と同様、BCCWJ コーパス中の共参照関係となる表現対のうち、文字列が完全一致しない366事例から50事例をランダムに抽出し、その事例の特徴を分析した。この結果を表13にまとめる<sup>6</sup>。

表13より、まず、教科書コーパスと比較して文字列完全一致しない共参照関係に代名詞が使われていることがわかる。これは、対象としたコーパス中の特に書籍レジスタ内に、人名や地名を指す代名詞、例えば、「彼」、「そこ」等の表現が多く出現したためだと考えられる。このように、BCCWJ コーパスには新聞ではほとんど出現しない、もしくは出現したとしても偏って出現する代名詞が多く含まれるため、今後アノテーションが完了し、研究分野で共有されることで、代名詞に対する照応解析の分析・評価のための有益な言語資源となることが予想される。

<sup>6</sup>1つの事例を複数の特徴に分類している場合も存在するため、分類結果の総数は50を越える点に注意されたい。

次に多かった特徴は先行詞とは異なる文字列の名詞句が共参照関係になる場合で、例えば、例(22)の「(バラナシの)施設」と「(マザーの)ところ」の関係、例(23)の「息子」と「サ우드」の関係がこれに該当する。また、この人物名「サ우드」は他の文脈では「魔術師」のような表現で記述されるため、そのような役割、肩書等を動的に獲得、理解しながら共参照の関係を同定する必要がある。

(22) 三ヶ月ほどをバラナシの施設で過ごしたのち、私はカルカッタへ移動した。するとカルカッタの宿で偶然、バラナシで仲良くしていたドイツ人のザビーナに出会ったのだった。「久しぶりね。あなたバラナシではマザーのところで働いていたでしょ。

(23) 王は息子の頭をなでて、  
「そうだね。わたしもみてみたいよ。きっと美しい宝石や杖なのだろう。けれどサ우드、よくおぼえておおき。...

また、例(22)の名詞句対「(バラナシの)施設」と「(マザーの)ところ」の共参照関係は複数の文節からなる表現間での共参照関係を捉える必要があるため、この名詞句の範囲を揺れなくアノテーションする基準や、さらに、複数の名詞句で書かれた内容の解釈を行った上で共参照解析を行うという、より複雑な処理が必要になると考えられる。

共参照関係のアノテーション時には、単語の切れ目の情報をアノテーション作業者に提示しなかったため、アノテーションされたタグに基づいて評価用データを作成する際にも問題が生じた。例えば、例(24)では、「我が国」と「日本」の間に共参照関係がアノテーションされているが、形態素解析の結果は「日本人」が最小の単位となっているため、誤って「我が国」と「日本人」に関係が付けられたデータとして扱ってしまっている。

(24) 近年、我が国における外国人の増加や諸外国との国際交流の進展により、日本語学習者は増加しており、海外で約235万人（平成15年国際交流基金調べ）、国内で13万5,146人（15年11月文化庁調べ）に上っています。

もちろん、作業者が付けた単語境界の情報を信頼し、それに基づいて作業結果を反映することも考えられるが、必ずしも作業者間で単語の切れ目に関する判断が一致するとは限らないため、例えば、単語境界をあらかじめ提示した上で、共参照関係をアノテーションを行うことも考えられるが、それでは本来正しく付けられるはずであった共参照関係が形態素解析の謝りのためにアノテーションされないという危険性も存在する。今後はWebページに記述されたテキスト等、形態素解析が必ずしも

高い精度で結果を出力されることが保証されないデータに対してアノテーションを行う需要がさらに高まると考えられるため、そのようなテキストに対しても単語境界を頑健に解析できる形態素解析器の出現が期待される。

また、照応詞もしくは先行詞が会話文の中に出現することで、解析が困難になる場合が存在する。例えば、例(25)の「女の子」と「私」は共参照関係となるが、この関係を捉えるためには、この括弧内の内容が「女の子」によって話されたものであることを特定する必要がある。

(25) 同じ宿の西洋人の女の子が、「あら、マザーテレサの施設ならここにもあるよ。私毎日行ってるから明日一緒に行かない?」と言ってくれたのだ。

この例では括弧の直前に出現した名詞句が発話者であることや、括弧の後に明示的に「言う」という表現が出現しているため、この括弧が会話文であり、その発話者の特定も容易であるが、文章のスタイルによっては必ずしも括弧が記述されない、「言う」のような明示的に発話であることを示す表現が存在しない、発話者が括弧の近傍に出現しない等の理由により、会話文とその発話者を特定することが困難な場合がある。このため、そのような場合にでも、頑健に発話者等の情報を特定できる処理が必要となる。

さらに、いわゆる節照応も問題となる。例えば、例(26)では「彼が聞き込みを押しつけた(こと)」と「そのこと」、例(27)では「そしてシスター達は私達ボランティアがどういう事をして一切文句を言わず、又お礼も言わなかった(こと)」と「それ」が照応関係となる。

(26) 彼に聞き込みを押しつけたのはステイシーだったが、いまはそのことを悔いていた。

(27) そしてシスター達は私達ボランティアがどういう事をして一切文句を言わず、又お礼も言わなかった。そしてそれは、そこに収容されている人達も同じだった。

節照応の関係をアノテーションする際の基準として節照応が出現した場合には、指し先の範囲特定の難しさの問題を回避するために、その範囲を特定せずに最右の文節を選択するという基準を採用している。しかし、これは便宜上こうやっているだけであり、アノテーションされた最右の文節をルートとする係り受け木全体が照応詞の指し先となる保証は無い。逆に、アノテーション作業者にその範囲の特定をまかせた場合には、作業者間で範囲特定の揺れが生じることが予測されるため、正確に節照応の関係をアノテーションを行うにはさらなる検討が必要となる。また、この関係を自動的に特定することに関しても何を手がかりに節照応として同定するかは自明ではない。まずは、今後節照応の関係にある事例を収集し、

その事例集合に共通して存在する特徴の分析が必要だと考えられる。このために、照応・共参照関係をアノテーションする際には、その再分類（例えば、節照応の関係にある、代名詞が照応詞となっている、等）も自動的にもしくは人手でアノテーションし、その再分類に基づいた分析や評価を行うべきだと考えられる。

## 6 おわりに

本稿では、省略関係と共参照関係がアノテーションされた世界史Bの教科書とBCCWJのコーパスの一部を対象に分析を行った。特に、述語項構造ならびに共参照関係を同定するツールであるKNPとSynChaを用い、現在利用可能なツールで解析できる事例と解析が困難である事例の特徴を分類し、その頻度を調査した。省略解析については、いくつかの観点から人手で解析に成功した事例や解析を誤った事例を分析し、省略解析技術の重要な構成要素となる選択制限・選択選好や主語のスコープについて独立してその性能を評価することで、省略解析手法の中で統合される個々の構成素の性能評価を行う方向性について述べた。選択制限・選択選好についてはKNPやSynChaに実装されたものを含め、複数の格を考慮した手法[10]等、さまざまな手法が提案されているため、それらの中で最も有効な手法を省略解析の文脈、つまり、実際の省略解析の問題へさらに選択制限の情報をアノテーションすることで各手法を評価していきたい。

また、省略解析の問題については、その関係の特定の困難さや定義の難しさのために起こるアノテーションの揺れの問題を避けて通ることができない。本稿で対象とした2種のコーパスに加え、NAISTテキストコーパスにも少なくないアノテーションの揺れや誤りが含まれるため、それらを容易に解消できる簡潔で、かつ、解析結果を利用する応用処理とも整合性の取れたアノテーション結果のクリーニングについても検討する必要がある。これについては、精密に仕様を確定する方法とは独立に、文献[23]で述べた複数名のアノテーション作業者を使った問題の定義が考えられる。これについては、アノテーション作業者の単純な多数決、もしくはクラウドソーシングで利用される作業者の信頼性も考慮した重み付きの多数決等、複数人の判定をどのように利用するかについてコンセンサスが得られる結果を得るために、どのような工夫をすべきかについて検討の余地がある。

また、共参照解析の分析では、代名詞の照応解析や、別称の知識獲得、節照応解析等、今後個別に考えていくべき課題について説明した。また、括弧表記の扱いや、会話文とその話者の特定、「ほとんど」や「一部」のような表現の意味的特性の継承は省略、共参照解析だけではなく、構文解析や情報抽出等、他のタスクの処理とも関

係があると考えられるため、分野横断的な統一的な扱いを共有することが重要であると考えられる。最後に、今後、括弧の解析や代名詞照応解析、節照応解析、本稿では扱わなかった外界照応解析は個別の問題として扱われ、その結果はある大学・研究機関等で独立に研究が進められることが予想されるが、それらを汲み上げ、統合するような共同研究基盤となる枠組み・ツールの開発が重要となると考えられる。これについては、他の研究分野からあがってくる共有すべき研究課題とともに今後どのように取り扱い、共有すべきかについて、本ワークショップで議論したいと考えている。

## 謝辞

KNPの使い方について詳しくご説明いただいた東京工業大学の笹野遼平助教に感謝いたします。

## 参考文献

- [1] Iida, R., Inui, K. and Matsumoto, Y.: Anaphora resolution by antecedent identification followed by anaphoricity determination, *ACM Transactions on Asian Language Information Processing (TALIP)*, Vol. 4, No. 4, pp. 417–434 (2005).
- [2] Iida, R., Inui, K. and Matsumoto, Y.: Zero-anaphora resolution by learning rich syntactic pattern features, *ACM Transactions on Asian Language Information Processing (TALIP)*, Vol. 6, No. 4, pp. 1–22 (2007).
- [3] Iida, R., Inui, K. and Matsumoto, Y.: Capturing Saliency with a Trainable Cache Model for Zero-anaphora Resolution, *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pp. 647–655 (2009).
- [4] Iida, R. and Poesio, M.: A Cross-Lingual ILP Solution to Zero Anaphora Resolution, *Proceedings of ACL-HLT 2011*, pp. 804–813 (2011).
- [5] Imamura, K., Saito, K. and Izumi, T.: Discriminative Approach to Predicate-Argument Structure Analysis with Zero-Anaphora Resolution, *Proceedings of ACL-IJCNLP*, pp. 85–88 (2009).
- [6] Miyao, Y. and Kawazoe, A.: University Entrance Examinations as a Benchmark Resource for NLP-based Problem Solving, *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pp. 1357–1365 (2013).
- [7] Nakaiwa, H. and Shirai, S.: Anaphora Resolution of Japanese Zero Pronouns with Deictic Reference,

- In Proceedings of the 16th COLING*, pp. 812–817 (1996).
- [8] Okumura, M. and Tamura, K.: Zero Pronoun Resolution in Japanese Discourse Based on Centering Theory, *In Proceedings of the 16th COLING*, pp. 871–876 (1996).
- [9] Taira, H., Fujita, S. and Nagata, M.: A Japanese Predicate Argument Structure Analysis using Decision Lists, *Proceedings of EMNLP*, pp. 523–532 (2008).
- [10] Van de Cruys, T.: A Non-negative Tensor Factorization Model for Selectional Preference Induction, *Proceedings of the Workshop on Geometrical Models of Natural Language Semantics*, pp. 83–90 (2009).
- [11] Yamazaki, K. M. M., Ogiso, T., Ogura, T. M. H., Kashino, W., Koiso, H., Yamaguchi, M., Tanaka, M. and Den, Y.: Balanced corpus of contemporary written Japanese, *Language Resources and Evaluation*, Vol. 48, No. 2, pp. 345–371 (2014).
- [12] 河原大輔, 黒橋禎夫, 橋田浩一: 「関係」タグ付きコーパスの作成, 言語処理学会第8回年次大会発表論文集, pp. 495–498 (2002).
- [13] 橋本力, 黒橋禎夫, 河原大輔, 新里圭司, 永田昌明: 構文・照応・評価情報つきブログコーパスの構築, 自然言語処理, Vol. 18, No. 2, pp. 175–201 (2011).
- [14] 光田航, 飯田龍, 徳永健伸: アノテーションとアノテーション作業者の信頼性推定, 言語処理学会第21回年次大会発表論文集 (2015).
- [15] 工藤拓, 市川宙, 中川哲治, 賀沢秀人: A joint inference of deep case analysis and zero subject generation for Japanese-to-English statistical machine translation, 言語処理学会第20回年次大会発表論文集, pp. 290–293 (2014).
- [16] 笹野遼平, 河原大輔, 黒橋禎夫, 奥村学: 構文・述語項構造解析システム KNP の解析の流れと特徴, 言語処理学会第19回年次大会発表論文集, pp. 110–113 (2013).
- [17] 笹野遼平, 黒橋禎夫: 自動獲得した名詞関係辞書に基づく共参照解析の高度化, 自然言語処理, Vol. 15, No. 5, pp. 99–118 (2008).
- [18] 笹野遼平, 黒橋禎夫: 大規模格フレームを用いた識別モデルに基づく日本語ゼロ照応解析, 情報処理学会論文誌, Vol. 52, No. 12, pp. 3328–3337 (2011).
- [19] 山本和英, 隅田英一郎: 決定木学習による日本語対話文の格要素省略補完, 自然言語処理, Vol. 6, No. 1, pp. 3–28 (1999).
- [20] 村田真樹, 長尾真: 用例や表層表現を用いた日本語文章中の指示詞・代名詞・ゼロ代名詞の指示対象の推定, 自然言語処理, Vol. 4, No. 1, pp. 87–109 (1997).
- [21] 萩行正嗣, 河原大輔, 黒橋禎夫: 外界照応および著者・読者表現を考慮した日本語ゼロ照応解析, 自然言語処理, Vol. 21, No. 3, pp. 563–600 (2014).
- [22] 萩行正嗣, 河原大輔, 黒橋禎夫: 多様な文書の書き始めに対する意味関係タグ付きコーパスの構築とその分析, 自然言語処理, Vol. 21, No. 2, pp. 213–247 (2014).
- [23] 飯田龍, 橋本力, 鳥澤健太郎, 黒橋禎夫, 乾健太郎, 宮尾祐介, 柴田知秀, 笹野遼平: 日本語書き言葉を対象とした人間の自然な省略検出の分析, 言語処理学会第21回年次大会発表論文集 (2015).
- [24] 飯田龍, 小町守, 井之上直也, 乾健太郎, 松本裕治: 述語項構造と照応関係のアノテーション: NAIST テキストコーパス構築の経験から, 自然言語処理, Vol. 17, No. 2, pp. 25–50 (2010).