

科学研究費助成事業 研究成果報告書

平成 28 年 6 月 1 日現在

機関番号：13901

研究種目：若手研究(B)

研究期間：2013～2015

課題番号：25730134

研究課題名(和文)音声言語アプリケーションのための漸進的係り受け解析技術の開発

研究課題名(英文)Development of incremental dependency parsing techniques for spoken language applications

研究代表者

大野 誠寛(Ohno, Tomohiro)

名古屋大学・情報基盤センター・助教

研究者番号：20402472

交付決定額(研究期間全体)：(直接経費) 3,100,000円

研究成果の概要(和文)：本研究では、同時通訳や字幕生成などの音声言語アプリケーションによって必要とされる新たな漸進的係り受け解析器を開発した。本解析器は、文節が入力されることに、既入力の文節間の依存関係、及び、係り先が未入力の文節に対しては既入力の文節との非依存性を同定する。また、係り先が未入力の文節が複数ある場合は、それらの係り先が同一か否かを同定する。次に、漸進的係り受け解析器の性能改善を目的に、人間の言語処理過程を表出した大規模データを構築し、人間による言語解析の能力や振舞の一端を明らかにした。さらに、話し言葉に対する係り受け解析精度の向上を目指して、係り受け解析と語順整序を同時に実行する手法を開発した。

研究成果の概要(英文)：In this research, we developed a new incremental dependency parser which is necessary for a spoken language application such as simultaneous interpretation and real-time captioning. Our parser identifies not only dependencies between inputted bunsetsus but also independencies between bunsetsus of which the modified bunsetsu has not been inputted and any other inputted bunsetsus, whenever a bunsetsu is inputted. In addition, if there are two or more bunsetsus of which the modified bunsetsu has not been inputted, our parser detects whether or not those bunsetsus depend on a same non-inputted bunsetsu. Next, to improve the performance of an incremental dependency parser, we constructed a large-scale dataset which expresses a process of human language processing, and revealed parts of human capacities and behaviors on incremental dependency parsing. Furthermore, we developed a method which can concurrently execute dependency parsing and word reordering.

研究分野：自然言語処理

キーワード：依存構造解析 漸進的処理 構文解析 話し言葉処理 字幕生成 リアルタイム処理 入力予測 語順整序

1. 研究開始当初の背景

同時通訳や字幕生成などの音声言語アプリケーションにおいて、キーワードベースの処理から脱却し、より豊かな処理を実現するためには、係り受け情報（構文情報）を利用することが不可欠である。これらのアプリケーションでは音声入力に追従した出力が求められるため、係り受け解析を利用するためには、音声入力と同時的に処理を進める解析技術が必要となる。

このような要請に答えるべく、これまでも、解析処理を漸進的に進めていく係り受け解析技術が開発されてきた[1~3]。しかし、いずれの研究も、従来の係り受け解析の枠組みと同様に、係り（係り元）と受け（係り先）の組を同定する問題として係り受け解析を捉え、その同定のための処理をどのように入力に追従して進めていくのが焦点となっている。いずれの手法も、入力が進むごとに解析の途中結果を更新しつつ保持しているが、入力途中の段階において、係り先が未だ入力されていない文節の係り受けに関する情報をどのように出力するのかについては、ほとんど検討されていない。あくまで、係りと受けの組をそれが同定できた段階で出力することを想定しており、その枠組みの中で精度や遅延時間を評価しているにすぎない。

このような従来の係り受け解析の枠組みを利用することを考えると、日本語の場合、係り先の文節はその係り元より後方に現れるという性質が強く存在するため、係りと受けの組の情報を取得できるタイミングは最速でも係り先文節の入力時間となる。すなわち、音声言語アプリケーションは、いつ入力されるとも分からない係り先文節が入力されるまで、何も情報を得られないまま処理を中断することになり、音声に追従した出力を実現することは難しくなる。

一方、応募者はこれまで、読みやすい字幕をリアルタイムに生成するための改行挿入技術の開発に取り組んできた。これらの開発では、従来の係り受け解析器を利用しているが、上述したように、その出力タイミングや出力内容は必ずしも利用しやすいものではなかった。この開発経験を通じて、より利用しやすく、より有益な係り受け解析の出力方式が他にあるのではないかと考えるに至った。

<参考文献>

- [1]加藤 他: 主辞情報付き文脈自由文法に基づく漸進的な依存構造解析, 信学論 (2003).
- [2]Johansson and Nugues: Incremental dependency parsing using online learning, In Proc. EMNLP-CoNLL2007 (2007).
- [3]Nivre: Algorithms for deterministic incremental dependency parsing, Computational Linguistics (2008).

2. 研究の目的

本研究では、同時通訳や字幕生成などの音声言語アプリケーションによって必要とされる係り受け解析の新たな出力方式を提案し、それを実現する解析器を開発することを当初の目的とした。その特徴は、係りと受けのペアの情報だけでなく、アプリケーションに有益なその他の係り受け情報を任意の時点で出力できる点にある。具体的には、以下に示す3つのサブゴールを設定し、徐々に高度な情報の出力が可能になるように漸進的係り受け解析器の開発を段階的に進めることとした。

- (1) 従来の係り受けペアに加えて、係り先が未入力の文節については入力済みの文節に係らないことを明示した係り受け構造（図1）を、任意の時点で出力可能な係り受け解析器（以下、解析器 a）を開発する。

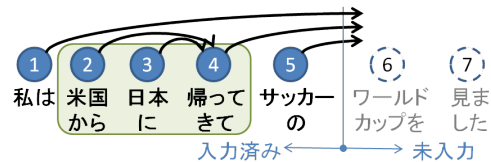


図1 解析器 a で出力する構造

- (2) 解析器 a を拡張して、係り先が未入力の文節が複数ある場合、それらの係り先が同じであることを明示した係り受け構造（図2）を、任意の時点で出力可能な係り受け解析器（以下、解析器 b）を開発する。

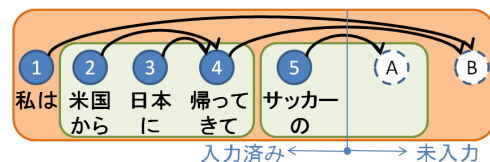


図2 解析器 b で出力する構造

- (3) 解析器 b を拡張して、係り先が未入力の文節に対して、その係り先が何であることを明示した係り受け構造（図3）を、任意の時点で出力可能な係り受け解析器（以下、解析器 c）を開発する。

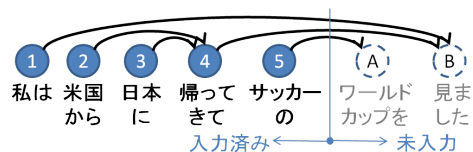


図3 解析器 c で出力する構造

3. 研究の方法

本研究は、平成25年度から平成27年度までの3年間で、音声言語アプリケーションの基盤技術として、漸進的係り受け解析技術を開発した。具体的には、上述した解析器 a から解析器 c までの開発及び検討を段階的に進

めた。以下では、各解析器の開発における研究方法について述べる。

(1)解析器 a の開発では、以下の ~ を順に実施した。

解析器 a の検討・実装：

まず、図 1 の構造の形式的な定義を行った。次に、従来の係り受け解析手法のうち、代表的な手法をいくつか取り上げ[4, 5]、本研究における漸進的係り受け解析に改良可能かどうか、また、その親和性について、定性的な検討を行った。その検討結果に基づき、各手法の解析途中結果から、任意の時点で図 1 の構造を出力できるように改良した。ただし、単に途中結果を出力するのではなく、入力済み文節には係らない確率を統計的に計算する機構を組み込むことにより、「入力済み文節には係らない」ことを同定する際の精度向上を図った。

解析器 a の評価実験：

まず、解析器 a の性能を測るため、図 1 の構造が文節入力ごとに出力されることを考慮した評価指標（解析精度）を新たに定義した。次に、係り受け情報付きの講演音声コーパスを用いて解析実験を実施し、解析器 a の解析精度を定量的に評価した。最後に、リアルタイム字幕生成のための逐次的な改行挿入手法[6]に対して、図 1 の構造から得られる構文的まとまりに関する情報を利用できるように改良を加え、その逐次的な改行挿入手法に解析器 a を適用することにより、改行挿入実験を実施した。実験データには、改行情報付き講演音声コーパスを用い、改行挿入の精度を評価した。

(2)解析器 b の開発では、以下の ~ を順に実施した。

解析器 b の検討・実装：

まず、図 2 の構造の形式的な定義を行った。次に、解析器 a を拡張し、図 2 の構造を任意の時点で出力可能な解析器 b を開発した。ここで、京都大学格フレームから得られる情報を素性として用いることにより、係り先が未入力である二つの文節の係り先が同一か否かを同定する際の精度向上を図った。格フレームを利用することにより、例えば、入力済みの二つの文節が格要素であり、両文節が一つの格フレームに含まれるならば、それらの文節が未だ入力されていない述語と係り受け関係をもつ可能性を見い出せるためである。

解析器 b の評価実験：

まず、解析器 b の性能を測るため、図 2 の構造が文節入力ごとに出力されることを考慮した評価指標（解析精度）を新たに定義した。次に、係り受け情報付きの講演音声コーパスを用いて解析実験を

実施し、解析器 b の解析精度を定量的に評価した。

(3)解析器 c の開発・検討では、以下の ~ を実施した。

解析器 c の検討：

まず、図 3 の構造の形式的な定義を行った。次に、図 3 の構造を任意の時点で出力可能な解析器 c を開発するため、解析器 b の性能と性質について検討を実施した。解析器 b にテキスト入力予測機能を単純に追加するだけでは、高い解析精度や入力予測精度を実現することは難しく、人間による漸進的な係り受け解析の能力や振舞を分析すること、また、語順を考慮した係り受け解析を実現することが先決であるとの結論を得た。そこで、解析器 c 自体の開発は保留し、下記に述べる ~ とを並行して進めた。

人間の言語処理過程を表出した大規模データを構築と分析：

まず、人間の言語理解過程における漸進性について指摘している先行研究について調査し、広く一般の文を対象として、その係り受け構造を漸進的に解析するという観点から、人間の言語処理過程を大規模に分析した研究はないことを確認した。その上で、構築するデータの設計方針を検討し定めた。次に、効率的なデータ構築を行うため、Web インタフェースを開発した。さらに、作業者を選定し、事前に Web インタフェースの使用テストを実施した。このテストにおいて、作業者には、設計方針やタグ付け仕様を確認・学習してもらうと同時に、Web インタフェースの使用感を確認してもらった。その後、作業者の意見に基づいて Web インタフェースの改修を図った上で、実際のデータ構築を実施した。構築したデータを用いて、人間による言語解析の能力や振舞に焦点を当てた分析を実施した。係り受け解析と語順整序の同時実行手法の開発：

まず、推敲支援や文生成などへの応用を目的とした語順整序に関する研究について調査した。次に、語順整序と係り受け解析を単に逐次的に実行するのではなく、同時的に実行する形での統合方法を考案し実装した。次に、実装した係り受け解析と語順整序の同時実行手法の性能を測るための実験方法を検討した。その結果、新聞記事中の文は読みやすい語順で書かれていることを前提に、文意は取れるものの読みにくい語順の文を新聞記事文から擬似的に作成し、それを実験のテストデータとした。なお、語順を機械的にランダムに変更しただけでは、母語話者が到底書きそうにない語順となる可能性があるため、人の判断を介在させて作成した。このようにして作成したテスト

データに対して、提案手法がどの程度、係り受け構造を同定できるか、また、元の新聞記事の語順を再現できるかを定量的に評価した。また、読みやすさの改善に関して主観的評価を実施した。

<参考文献>

- [4]内元 他: 後方文脈を考慮した係り受けモデル, 自然言語処理 (2000).
- [5]工藤, 松本: チャンキングの段階適用による日本語係り受け解析, 情処学論 (2002).
- [6]大野 他: 講演のリアルタイム字幕生成のための逐次的な改行挿入, 電学論 (2013).

4. 研究成果

本研究では主に、平成 25 年度～平成 27 年度の 3 年間で、(1)解析器 a の開発、(2)解析器 b の開発、(3)解析器 c の検討、(4)人間の言語処理過程を表出した大規模データの構築と分析、(5)係り受け解析と語順整序の同時実行手法の開発、を実施し、以下に示す研究成果を得ることができた。

(1) 解析器 a の開発:

音声言語アプリケーションが文節間の依存関係に関する情報をできる限り早期に利用できることが望ましいという観点から、漸進的な係り受け解析が生成する新たな係り受け構造を提案し、その構造を出力可能な解析器 a を開発した。本解析器は、文節が入力されるごとに解析を実行し、既入力文節間の依存関係、ならびに、係り先が未入力文節に対しては既入力文節との非依存性を同定し、それを明示した係り受け構造を生成する。係り先が未だ入力されていないという情報を解析器が提示することにより、上位層のアプリケーションは文節列の構文的なまとまりの成否をより考慮した処理が可能となる。係り受け解析実験により、本解析器が、文を単位とする解析手法と比べて精度を大きく低下させることなく、漸進的な係り受け解析を実現できることを示した。また、リアルタイム字幕生成における改行位置の同定に本解析器を応用し、本研究で提案した係り受け構造の有効性を確認した。

以上の研究成果は、言語処理学会第 20 回年次大会で発表した後、電子情報通信学会論文誌、国際会議 IWPT2013 (The 13th International Conference on Parsing Technologies) に採録・採択されており、国内外で高い評価を得ることができた。

(2) 解析器 b の開発:

解析器 a を拡張し、未入力文節との構文的関係を明示する漸進的な係り受け解析器 b を開発した。本解析器は、文節が入力されるたびに解析を実行し、係り先

が入力されていない文節に対しては、係り先が未入力であることを同定する。さらに、係り先が未入力である文節が複数あるときは、それらの係り先が同一か否かを同定する。日本語講演データを用いた解析実験の結果、係り先が正しく同定できているかの判定において正解率 71.22%を達成し、未入力文節との構文的関係が明示された係り受け構造を精度よく生成できることを確認した。

以上の研究成果は、言語処理学会第 20 回年次大会で発表した。

(3) 解析器 c の検討:

解析器 b を拡張することにより、係り先が未入力文節に対して、その係り先がどのような文節であるかを明示した係り受け構造を任意の時点で出力可能な係り受け解析器 c の開発のための検討を進めた。その結果、高い解析精度や入力予測精度を実現するためには、人間による漸進的な係り受け解析の能力や振舞を分析すること、また、語順を考慮した係り受け解析を実現することが先決であるとの結論を得た。そのため、解析器 c の開発は今後の課題とし、以下に示す項目(4)と(5)を先行して実施した。

(4) 人間の言語処理過程を表出した大規模データの構築と分析:

人間による漸進的な係り受け解析結果を定量的に分析し、その能力や振舞を明らかにすることができれば、漸進的係り受け解析器の性能を向上させるための知見が得られる可能性がある。そこで漸進的解析器の性能改善を目的に、人間の言語処理過程を表出した大規模データを構築した。本データの特徴は、京都大学テキストコーパスに収録された新聞記事中の 2,502 文を対象として、文節が文頭から順に 1 つ提示されるたびに、それまでに提示された文節列に対して、人間が係り受け解析と入力予測を施している点にある。なお、本データ構築では、事前に Web インタフェースを開発し、効率的な作業を実施している。

さらに、構築したデータを分析し、人間による言語解析の能力や振舞の一端を明らかにした。具体的には、人間にとって、係り先が未入力文節であることを単に示すこと(図 1 の構造の解析)はそれほど難しいタスクではないが、その未入力文節との間の係り受け関係を示すこと(図 2 の構造の解析)が必要となると途端に難しいタスクとなることがわかった。さらに、未入力文節の係り先文節の文字列を予測することは、人間にとって難しいタスクであるものの、その一方で、1 割程度の未入力文節についてはその文字列(の一部)を予測できており、全く予測

できないわけではないことがわかった。

以上の研究成果は、言語処理学会第22回年次大会で発表した。

(5) 係り受け解析と語順整序の同時実行手法の開発：

入力文が読みにくい語順である場合、係り受け解析の精度は低下する傾向にあり、話し言葉に対する係り受け解析精度が低くなる一因となっていた。そこで、話し言葉に対する係り受け解析精度の向上を図ることを目標として、語順と係り受けは相互に関連していることを考慮し、係り受け解析と語順整序を同時に実行する手法を開発した。本手法は、係り受け構造が付与されていない文を入力とし、係り受け解析と語順整序を同時に行う。係り受けと語順の適切さを同時に考慮することにより、読みやすい語順に精度よく整えることができる。

評価実験では、文法的には間違っていないものの読みにくい語順を持つ文を新聞記事から擬似的に552文作成し、それらに対して語順整序を実行した。元の新聞記事文の語順を正解としたときの、語順整序結果との一致を測定した結果、本手法による語順整序結果は、比較のために設定した2つのベースラインと比べ高い正解率(83.82%)を達成しており、本手法の有効性を確認した。一方、本手法と両ベースラインの係り受け解析精度を比較評価した結果、本手法の係り受け単位正解率は、両ベースラインと比べて低かった。一方、本手法の文単位正解率は、両ベースラインよりも上回った。本手法は、両ベースラインと比べて、文全体の係り受け構造の同定に失敗する場合は、その文において、より多くの係り受け関係の同定を誤ることになるが、その一方で、より多くの文において、文全体の係り受け構造を正しく同定できるという傾向をもつといえる。

以上の研究成果は、電子情報通信学会論文誌、国際会議 Coling2014 (The 25th International Conference on Computational Linguistics) 及び、国際会議 ENLG2015 (The 15th European Workshop on Natural Language Generation) に採録・採択されており、国内外で高い評価を得ることができた。

5. 主な発表論文等

(研究代表者、研究分担者及び連携研究者には下線)

〔雑誌論文〕(計2件)

大野 誠寛, 吉田 和史, 加藤 芳秀, 松原 茂樹: 係り受け解析との同時実行に基づく日本語文の語順整序, 電子情報通信学会論文誌, Vol. J99-D, No. 2, pp. 201-213,

Feb. (2016).

大野 誠寛, 松原 茂樹: 文節間の依存・非依存を同定する漸進的係り受け解析, 電子情報通信学会論文誌, Vol. J98-D, No. 4, pp. 709-718, Apr. (2015).

〔学会発表〕(計7件)

大野 誠寛, 松原 茂樹: 漸進的係り受け解析の出力構造 - 人間の文解析過程のアノテーション -, 言語処理学会第22回年次大会発表論文集, pp. 709-712, 東北大学・川内北キャンパス, Mar. (2016).

Tomohiro Ohno, Kazushi Yoshida, Yoshihide Kato, Shigeki Matsubara: Japanese Word Reordering Executed Concurrently with Dependency Parsing and Its Evaluation, *Proceedings of the 15th European Workshop on Natural Language Generation (ENLG2015)*, pp. 61-65, Brighton, UK, Sep. (2015).

吉田 和史, 大野 誠寛, 加藤 芳秀, 松原 茂樹: 係り受け解析との統合に基づく日本語文の語順整序, 情報処理学会研究報告, Vol. 2015-NL-220, No. 13, pp. 1-10, 九州大学医学部百年講堂, Jan. (2015).

Kazushi Yoshida, Tomohiro Ohno, Yoshihide Kato, Shigeki Matsubara: Japanese Word Reordering Integrated with Dependency Parsing, *Proceedings of the 25th International Conference on Computational Linguistics (COLING2014)*, pp. 1186-1196, Dublin, Ireland, Aug. (2014).

村田 匡輝, 大野 誠寛, 松原 茂樹: 未入力文節との構文的関係を考慮する漸進的な係り受け解析, 言語処理学会第20回年次大会発表論文集, pp. 193-196, 北海道大学, Mar. (2014).

吉田 和史, 大野 誠寛, 加藤 芳秀, 松原 茂樹: 係り受け解析を伴った日本語文の語順整序, 言語処理学会第20回年次大会発表論文集, pp. 701-704, 北海道大学, Mar. (2014).

Tomohiro Ohno, Shigeki Matsubara: Dependency Structure for Incremental Parsing of Japanese and its Application, *Proceedings of the 13th International Conference on Parsing Technologies (IWPT2013)*, pp. 91-97, Nara, Japan, Nov. (2013).

〔その他〕

ホームページ等

<http://slp.itc.nagoya-u.ac.jp/~ohno/kaken/incre-dep-parsing/index.html>

6. 研究組織

(1) 研究代表者

大野 誠寛 (OHNO, Tomohiro)

名古屋大学・情報基盤センター・助教
研究者番号： 20402472