

# ニュース記事に対する談話構造と興味度のアノテーション ～ニュース対話システムのパーソナライズに向けて～

高津 弘明<sup>1</sup> 安藤 涼太<sup>2</sup> 本田 裕<sup>3</sup> 松山 洋一<sup>1</sup> 小林 哲則<sup>1</sup>

<sup>1</sup> 早稲田大学

<sup>2</sup> 内外切抜通信社

<sup>3</sup> 本田技研工業

takatsu@pcl.cs.waseda.ac.jp ando@naigaipc.co.jp hiroschi\_01\_honda@jp.honda  
matsuyama@pcl.cs.waseda.ac.jp koba@waseda.jp

## 1 はじめに

談話構造に基づく文章の首尾一貫性を保持しつつ、ユーザごとにパーソナライズした情報を提供する音声対話システムの開発を目的として、ニュース記事に対して談話構造と文単位のユーザの興味度を付与したデータセットを構築した。

インターネットの普及により、日々膨大なデジタルコンテンツがウェブ上に蓄積され続けている。さらに、近年ではパソコンやスマートフォンなどが普及したことで、ウェブ上の情報に容易にアクセスできるようになり、人々のメディアへの接触時間は年々増加している [1]。一方で「世の中の情報量が多すぎる」という意識も高まってきており、限られた時間の中でより効率的に情報を消費することへのニーズが高まっている。我々は、このような社会背景を踏まえ、人々が日々消費する情報の中でも特にニュースに焦点を当て、効率的にニュース記事の内容を伝達する音声対話システム（以下、ニュース対話システム）[2]の開発を進めている。

ニュース記事の内容を効率的に分かりやすく伝えるうえで重要な観点として、ほしい情報をどれだけ提示でき、いらぬ情報をどれだけ除外できるか（情報伝達効率）という点と、文章の構造に基づき意味的に整合性のある文のつながりを保持したまま情報を伝えられるか（首尾一貫性）という点がある。ニュース対話システムは、要約に基づく主計画と補足説明のための副計画から構成されるシナリオに基づいて会話を進行させる。ユーザが受け身で聴いている限りにおいては主計画の内容を伝えることになる。主計画として、重要度に基づきユーザごとに共通の要約を提示するよりも、興味度に基づきユーザごとにパーソナライズした要約を提示する方が高い情報伝達効率を実現できることを確認し

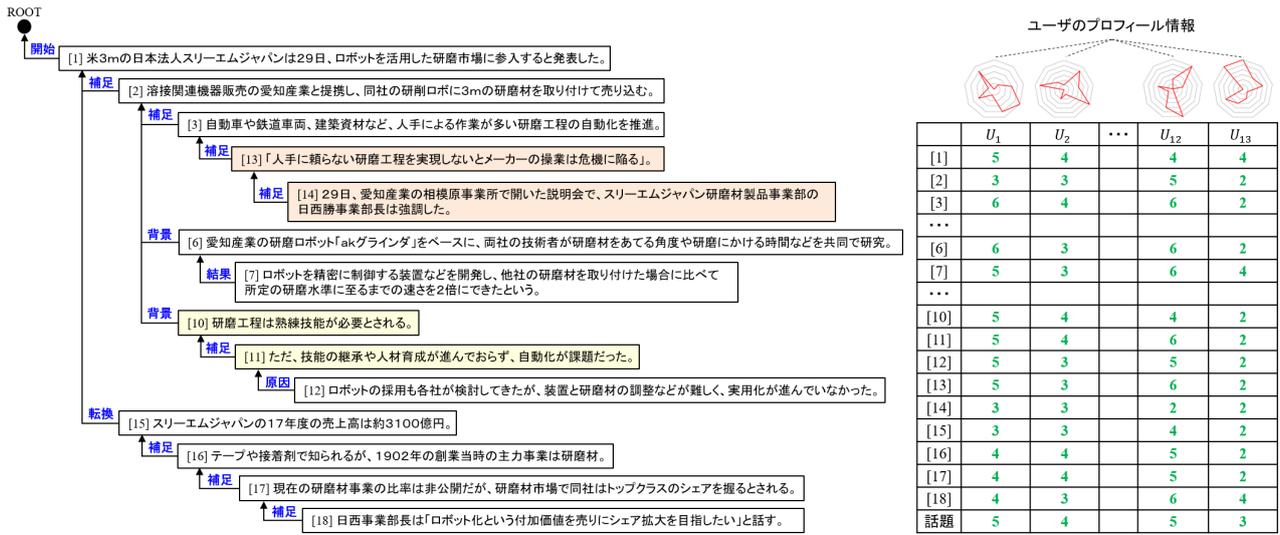
た [3]。一方で、ニュース内容を適切に理解させるためには話の内容が首尾一貫している必要がある。

そこで、文章の首尾一貫性を考慮しつつ、ユーザの興味がある情報を提供する音声対話システムを開発するために、日本語のニュース記事に対して図1のような談話構造とユーザの興味度を付与したデータセットを構築した。談話構造に関しては、文を談話単位とし、「係り先」「談話関係」「チャンク」のアノテーションを行った。例を図1(a)に示す。係り先は、ルートノードから親ノード文までの情報が対象の文を理解するうえで必要最小限の情報になるように付与したものである。談話関係は、子ノード文が親ノード文に対してどのような意味的關係にあるかを分類したものである。チャンクは、親ノード文の内容を正しく理解するために欠かせない情報が子ノード文に書かれている場合、これらをまとめて提示すべきであることを表したものである。合計1200記事に対して専門家2人が談話構造を付与した。興味度に関しては、クラウドソーシングを用いて、様々な属性の被験者に、ニュース記事の話題と各文の内容について興味度を6段階で回答させた。例を図1(b)に示す。談話構造を付与した1200記事に対して、1人あたり6記事、1記事あたり11人以上となるように配分し、合計2507人分のプロフィールと興味度のデータを収集した。

本稿の構成は次の通りである。2章で関連研究について述べる。3章で談話構造データセットについて説明する。4章で興味度データセットについて説明する。5章で今後の展望について述べる。

## 2 関連研究

談話構造解析は、文章を構成する文や節の間に成り立つ関係を解析する自然言語処理の基本的なタスクである。談話構造解析の結果は、文書要約 [4] や



(a) 係り先・談話関係・チャンク(強いチャンク・弱いチャンク)のアノテーション例 (b) 興味度のアノテーション例

図 1 アノテーションの例

質問応答 [5], 機械翻訳 [6], 評判分析 [7] などの下流タスクのアプリケーションで用いられる [8].

談話構造解析のための代表的なデータセットとして, RST Discourse Treebank [9], Discourse Graphbank [10], Penn Discourse Treebank [11] がある. RST Discourse Treebank は, 修辞構造理論 (Rhetorical Structure Theory) [12] に基づいて構築されたデータセットである. 修辞構造理論では, Elementary Discourse Unit (EDU) と呼ばれる節を最小の談話単位として, 隣接する EDU 同士を談話関係で結合し, より大きな談話単位を形成する. さらに, その談話単位同士も談話関係によって統合され, 最終的に修辞構造木が形成される. Discourse Graphbank は, グラフ構造を採用し, 1つの談話単位から複数の談話単位に談話関係が付与されたデータセットである. Penn Discourse Treebank は, 木構造やグラフ構造を仮定せず, 談話単位間の関係を 2 項関係で表現したデータセットである. 隣接する談話単位の間接続表現が存在する場合は, その接続表現を手がかりとして談話関係を分類する. 談話単位の間接続表現が存在しない場合は, 接続表現が挿入可能であるかを判断し, 談話関係を分類する. これらは英語のニュース記事に基づいて作成された談話構造のデータセットであるが, 英語以外にも, ドイツ語 [13], スペイン語 [14], ポルトガル語 [15], 中国語 [16] など様々な言語の談話構造のデータセットが構築されている. 6つの言語の対訳テキストに対して談話構造を付与したデータセットも構築されている [17].

日本語の文書に対して談話構造をアノテ

ションする研究も存在する [18, 19, 20]. 金子らは, BCCWJ [21] の文を分解して得られたセグメントに対して, 分節談話表示理論 (Segmented Discourse Representation Theory) [22] に基づいて時間関係と因果関係をアノテーションした [18]. 河原らは, 様々なドメインのウェブ文書の書き始め 3 文を対象にクラウドソーシングで談話関係のアノテーションを行う方法を提案し, 大規模な文書への談話関係のアノテーションを短時間でできることを示した [19]. さらに, 岸本らは, 河原らのアノテーション基準に言語テストを追加するなどの改良を加えることでアノテーションの品質が向上することを確認した [20].

文書要約 [23, 24, 25, 26] のようなタスクへの応用を考えた場合, 修辞構造木のような句構造よりも談話単位間の親子関係を直接表せる依存構造の方が望ましい. そこで, 修辞構造木を談話依存構造木へ変換する方法が提案された [23, 27]. 一方で, 変換アルゴリズムによって生成される談話依存構造木の性質が異なることから [28], Yang らは, 科学論文のアブストラクトに対して人手で EDU 間の依存構造と談話関係をアノテーションする方法を提案し, データセットとして SciDTB を構築した [29].

本研究では, 日本語のニュース記事に対して専門家が文単位の談話依存構造をアノテーションし, クラウドソーシングにより, ユーザのプロフィールおよびニュース記事の話題と文に対する興味度を収集することで, 文章の一貫性を考慮しながらパーソナライズした情報を伝達する要約・対話システムの構築に活用できるデータセットを作成した.

### 3 談話構造データセット

15文から25文の日本語のニュース記事1200個に対して、ウェブニュースのクリッピングの専門家2人が「係り先」「談話関係」「チャンク」のアノテーションを行った。ジャンルの内訳は「スポーツ」「テクノロジー」「経済・政治」「国際」「社会」「地域」であり、なるべく話題が重複しないように人手で各ジャンル200記事ずつ選択した。アノテーション作業は係り先、談話関係、チャンクの順番で行った。談話単位は文であり、文とは地の文に出現する句点“。”で区切られる文字列を表す。2人のアノテータは、半分ずつ作業を分担し、自分の作業を終えた後は他方のアノテータのアノテーション結果のチェックにまわり、疑問点が見つかった場合は両者の合意が得られるまで議論を重ねた。

#### 3.1 係り先のアノテーション

文 $j$ を文 $i$ の係り先として指定できる条件は以下の通りである。

- 原文において文 $j$ が文 $i$ よりも前に出現している。
- 木構造に従ってルートノードから順番に読み進め、文 $j$ の後に文 $i$ を読んだ際に、話の流れが自然である。
- ルートノードから文 $j$ までの情報が文 $i$ を理解するための必要最小限の情報になっている。
- 文 $i$ から読み始めることが可能である場合、文 $i$ の係り先はルートノードとする。

#### 3.2 談話関係のアノテーション

談話関係の定義を表1に示す。談話関係の定義に際して、日本語を対象とした先行研究の談話関係の分類[30, 31, 32, 33]を参考にしつつ、ニュース対話システムへの応用を考えた際の必要十分な粒度であることおよび、実際にニュース記事から構築された談話依存構造木において頻繁に観測される文間関係であることに留意した。アノテーションの判断は、談話関係の定義と対話基準の両方の条件に合致しているかを確認しながら行った。対話基準とは、談話関係に基づく応答が自然かどうかに基づく判断であり、例えば「原因」について判断する際、親ノード文の内容を伝えた後に「なんで？」と質問が来たときに、子ノード文の内容を回答として提示しても不自然でないことを確認する。なお、1つの文間に対して複数の談話関係が付与されることもある。

表1 談話関係の定義 (括弧内の説明は対話基準を表す)

談話関係	アノテーション基準
開始	文の係り先がルートノードである。
結果	子ノード文が親ノード文の結果である。(親ノード文に対する「それでどうなったの？」等の結果を問う質問の返答として子ノード文を提示することが適切である。)
原因	子ノード文が親ノード文の原因である。出来事の原因、筆者の主張の根拠、理由などが含まれる。(親ノード文に対する「なんで？」等の原因を問う質問の返答として子ノード文を提示することが適切である。)
背景	親ノード文で事実や出来事が述べられていて、子ノード文でそれらの背景となる事柄や前提が述べられている。(親ノード文に対する「どんな背景があるの？」等の背景を問う質問の返答として子ノード文を提示することが適切である。)
呼応	親ノード文の問いかけに対して、子ノード文で答えが示されている。問題や課題とそれに対する対策や検討なども含む。(親ノード文をユーザの疑問とみなしたとき、その返答として子ノード文を提示することが適切である。)
対比	親ノード文と子ノード文が対比関係にある。
転換	親ノード文から子ノード文にかけて話題が変わっている。サブトピックレベルの変化も含む。
例示	親ノード文で述べられた事柄の具体例が子ノード文で提示されている。(親ノード文に対する「例えば？」等の具体例を求める質問の返答として子ノード文を提示することが適切である。)
結論	子ノード文が親ノード文までの話のまとめや結論になっている。(親ノード文に対する「つまり？」等の結論を求める質問の返答として子ノード文を提示することが適切である。)
補足	親ノード文で述べられていた事柄についての詳細や補足が子ノード文で述べられている。広義には上記談話関係に含まれない補足関係にある文間関係に対して「補足」を付与する。(親ノード文に対する「詳しく教えて」等の補足要求の返答として子ノード文を提示することが適切である。)

#### 3.3 チャンクのアノテーション

親ノード文を提示する際、子ノード文も一緒に提示することが望まれる場合、これらをチャンクとする。特に親ノード文の内容を理解するうえで必須となる情報が子ノード文に書かれている場合、これらを「強いチャンク」とする。例えば図1のような、親ノード文に発言、子ノード文に発言者の情報が書かれている場合や、複数文に渡る手順説明などは強いチャンクとなる。また、子ノード文の情報は必ずしも必須ではないが、親ノード文の内容に対する偏った理解を防ぐために子ノード文の情報が重要である場合、これらを「弱いチャンク」とする。例えば図1のような、現状求められているものが親ノード文で述べられており、その実現が容易でないことを子ノード文で説明している場合や、主題に関わる2つの国の情勢について、一方の説明が親ノード文に書かれており、他方の説明が子ノード文に書かれている場合などは弱いチャンクとなる。

#### 3.4 データセットの統計

記事ごとのルートノードから葉ノードまでの文数の最大値(木の深さ)の分布を図2に示す。1記事

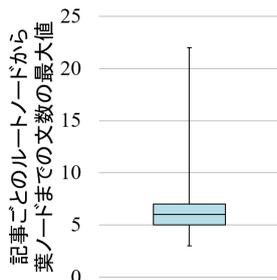


図2 木の深さの分布

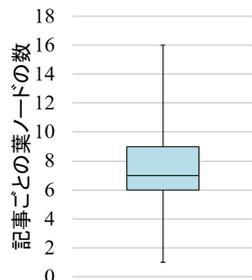


図3 木の幅の分布

あたり木の深さは平均で 6.5 文程度であった。記事ごとの葉ノードの数（木の幅）の分布を図3に示す。1 記事あたり木の幅は平均で 7.5 文程度であった。

データセットにおける談話関係の出現頻度は、「開始」が 1221 個、「結果」が 400 個、「原因」が 691 個、「背景」が 1343 個、「呼応」が 851 個、「対比」が 638 個、「転換」が 220 個、「例示」が 709 個、「結論」が 920 個、「補足」が 14609 個であった。

強いチャンクの数 は 231 個、弱いチャンクの数 は 699 個であった。また、1 チャンクあたりの平均文数は 2.15 文であった。

## 4 興味度データセット

クラウドソーシングを用いて被験者を募集し、以下のプロフィールアンケートと興味度アンケートに回答させた。ニュース記事には談話構造データセットで使ったものと共通の 1200 記事を使用した。1 人あたり 6 ジャンル 1 記事ずつ、1 記事あたり 11 人以上が回答するように配分し、2507 人分のアンケート結果を収集した。

### 4.1 プロフィールアンケート

以下の 14 項目について回答させた。性別、年齢、住んでいる地域、職種、業種、ニュースを見る頻度、ニュースをよくチェックする時間帯、映像・音声・文字のうちニュースへの接触方法として多いものはどれか、ニュースを知る手段、ニュースを読む際に使用している新聞やウェブサイト・アプリ、有料でニュースを読んでいるか、普段積極的に読む・見る・聞くニュースのジャンル、ニュースのジャンルに対する興味程度、趣味。

### 4.2 興味度アンケート

ニュース記事の本文を読ませ、内容を理解させた後、記事全体の話題と各文の内容について興味度を次の 6 段階で回答させた。6:とても興味がある、5:

■ まったく興味がない ■ 興味がない ■ どちらかという興味がない  
■ どちらかという興味がある ■ 興味がある ■ とても興味がある



(a)記事の話題に対する興味度の割合 (b)文の内容に対する興味度の割合

図4 アンケートで得られた興味度の割合

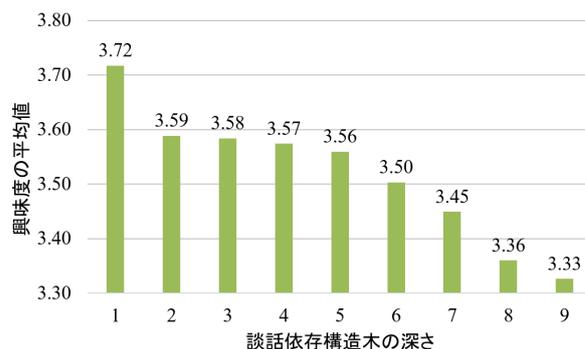


図5 談話依存構造木の深さと興味度の平均値

興味がある、4:どちらかといえば興味がある、3:どちらかといえば興味がない、2:興味がない、1:まったく興味がない。

### 4.3 データセットの統計

合計で 15042 記事 268509 文に対する回答を得た。記事の話題に対する興味度の割合および文の内容に対する興味度の割合を図4に示す。

談話依存構造木の深さごとに文の興味度の平均値を計算した結果を図5に示す。この結果から、多くの人の興味を引く内容は木の浅いところにあり、人によって興味が分かれる内容は木の深いところにあることが分かった。

## 5 おわりに

ニュース記事に対して談話構造および、ユーザのプロフィールと記事の話題・文に対するユーザの興味度を付与したデータセットを構築した。

今後は、本データセットを用いて、談話構造解析 [34] や、談話構造を制約とした抽出型要約のパーソナライズ [35]、談話関係に基づく質問応答や対話制御などについて検討する。

**謝辞** 本研究は、JST START (JPMJST1912) の支援を受けたものである。

## 参考文献

- [1] 博報堂 D Y メディアパートナーズ (編). 広告ビジネスに関わる人のメディアガイド 2020: メディア環境のこれからと今. 宣伝会議, pp. 20–40, 2020.
- [2] 高津弘明, 福岡維新, 藤江真也, 林良彦, 小林哲則. 意図性の異なる多様な情報行動を可能とする音声対話システム. 人工知能学会論文誌, Vol. 33, No. 1, pp. 1–24, 2018.
- [3] Hiroaki Takatsu, Mayu Okuda, Yoichi Matsuyama, Hiroshi Honda, Shinya Fujie, and Tetsunori Kobayashi. Personalized extractive summarization for a news dialogue system. In *Proceedings of the 8th IEEE Spoken Language Technology Workshop*, 2021.
- [4] Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu. Discourse-aware neural extractive text summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5021–5031, 2020.
- [5] Peter Jansen, Mihai Surdeanu, and Peter Clark. Discourse complements lexical semantics for non-factoid answer reranking. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pp. 977–986, 2014.
- [6] Junxuan Chen, Xiang Li, Jiarui Zhang, Chulun Zhou, Jianwei Cui, Bin Wang, and Jinsong Su. Modeling discourse structure for document-level neural machine translation. In *Proceedings of the 1st Workshop on Automatic Simultaneous Translation*, pp. 30–36, 2020.
- [7] Baptiste Chardon, Farah Benamara, Yannick Mathieu, Vladimir Popescu, and Nicholas Asher. Measuring the effect of discourse structure on sentiment analysis. In *Proceedings of the 14th international conference on Computational Linguistics and Intelligent Text Processing*, pp. 25–37, 2013.
- [8] Shafiq Joty, Giuseppe Carenini, Raymond T. Ng, and Gabriel Murray. Discourse analysis and its applications. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pp. 12–17, 2019.
- [9] Lynn Carlson, Daniel Marcu, and Mary Ellen Okurovsky. Building a discourse-tagged corpus in the framework of rhetorical structure theory. In *Proceedings of the 2nd SIGdial Workshop on Discourse and Dialogue*, pp. 1–10, 2001.
- [10] Florian Wolf and Edward Gibson. Representing discourse coherence: A corpus-based study. *Computational Linguistics*, Vol. 31, No. 2, pp. 249–287, 2005.
- [11] Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Miltsakaki, Livio Robaldo, Aravind Joshi, and Bonnie Webber. The penn discourse treebank 2.0. In *Proceedings of the 6th International Conference on Language Resources and Evaluation*, pp. 2961–2968, 2008.
- [12] William C. Mann and Sandra A. Thompson. Rhetorical structure theory: Toward a functional theory of text organization. *Text*, Vol. 8, No. 3, pp. 243–281, 1988.
- [13] Manfred Stede. The potsdam commentary corpus. In *Proceedings of the Workshop on Discourse Annotation*, pp. 96–102, 2004.
- [14] Iria da Cunha, Juan-Manuel Torres-Moreno, and Gerardo Sierra. On the development of the RST Spanish treebank. In *Proceedings of the 5th Linguistic Annotation Workshop*, pp. 1–10, 2011.
- [15] Thiago Alexandre Salgueiro Pardo, Maria das Graças Volpe Nunes, and Lucia Helena Machado Rino. DiZer: An automatic discourse analyzer for Brazilian Portuguese. In *Proceedings of the Brazilian Symposium on Artificial Intelligence*, pp. 224–234, 2004.
- [16] Lanjun Zhou, Binyang Li, Zhongyu Wei, and Kam-Fai Wong. The CUHK discourse treebank for Chinese: Annotating explicit discourse connectives for the Chinese treebank. In *Proceedings of the 9th International Conference on Language Resources and Evaluation*, pp. 942–949, 2014.
- [17] Deniz Zeyrek, Amalia Mendes, Yulia Grishina, Murathan Kurfali, Samuel Gibbon, and Maciej Ogrodniczuk. TED multilingual discourse bank (TED-MDB): A parallel corpus annotated in the PDTB style. *Language Resources and Evaluation*, Vol. 54, pp. 587–613, 2020.
- [18] Kimi Kaneko and Daisuke Bekki. Building a Japanese corpus of temporal-causal-discourse structures based on SDRT for extracting causal relations. In *Proceedings of the EACL 2014 Workshop on Computational Approaches to Causality in Language*, pp. 33–39, 2014.
- [19] Daisuke Kawahara, Yuichiro Machida, Tomohide Shibata, Sadao Kurohashi, Hayato Kobayashi, and Manabu Sassano. Rapid development of a corpus with discourse annotations using two-stage crowdsourcing. In *Proceedings of the 25th International Conference on Computational Linguistics*, pp. 269–278, 2014.
- [20] Yudai Kishimoto, Shinnosuke Sawada, Yugo Murawaki, Daisuke Kawahara, and Sadao Kurohashi. Improving crowdsourcing-based annotation of Japanese discourse relations. In *Proceedings of the 11th International Conference on Language Resources and Evaluation*, pp. 4044–4048, 2018.
- [21] Kikuo Maekawa, Makoto Yamazaki, Toshinobu Ogiso, Takehiko Maruyama, Hideki Ogura, Wakako Kashino, Hanae Koiso, Masaya Yamaguchi, Makiro Tanaka, and Yasuharu Den. Balanced corpus of contemporary written Japanese. *Language Resources and Evaluation*, Vol. 48, No. 2, pp. 345–371, 2014.
- [22] Nicholas Asher and Alex Lascarides. *Logics of conversation: Studies in natural language processing*. Cambridge University Press, 2003.
- [23] Tsutomu Hirao, Yasuhisa Yoshida, Masaaki Nishino, Norihito Yasuda, and Masaaki Nagata. Single-document summarization as a tree knapsack problem. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1515–1520, 2013.
- [24] Yasuhisa Yoshida, Jun Suzuki, Tsutomu Hirao, and Masaaki Nagata. Dependency-based discourse parser for single-document summarization. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pp. 1834–1839, 2014.
- [25] Yuta Kikuchi, Tsutomu Hirao, Hiroya Takamura, Manabu Okumura, and Masaaki Nagata. Single document summarization based on nested tree structure. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pp. 315–320, 2014.
- [26] Tsutomu Hirao, Masaaki Nishino, Yasuhisa Yoshida, Jun Suzuki, Norihito Yasuda, and Masaaki Nagata. Summarizing a document by trimming the discourse tree. *IEEE/ACM Transactions on Audio, Speech and Language Processing*, Vol. 23, No. 11, pp. 2081–2092, 2015.
- [27] Sujian Li, Liang Wang, Ziqiang Cao, and Wenjie Li. Text-level discourse dependency parsing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pp. 25–35, 2014.
- [28] Katsuhiko Hayashi, Tsutomu Hirao, and Masaaki Nagata. Empirical comparison of dependency conversions for RST discourse trees. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pp. 128–136, 2016.
- [29] An Yang and Sujian Li. SciDTB: Discourse dependency treebank for scientific abstracts. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 444–449, 2018.
- [30] 黒橋禎夫, 長尾眞. 表層表現中の情報に基づく文章構造の自動抽出. 自然言語処理, Vol. 1, No. 1, pp. 3–20, 1994.
- [31] 梅澤俊之, 原田実. センタリング理論と対象知識に基づく談話構造解析システム DIA. 自然言語処理, Vol. 18, No. 1, pp. 31–56, 2011.
- [32] 山本和英, 齋藤真実. 用例利用型による文間接続関係の同定. 自然言語処理, Vol. 15, No. 3, pp. 21–51, 2008.
- [33] 横山憲司, 難波英嗣, 奥村学. Support Vector Machine を用いた談話構造解析. 情報処理学会研究報告, Vol. 2003, No. 23(2002-NL-154), pp. 193–200, 2003.
- [34] 高津弘明, 安藤涼太, 松山洋一, 小林哲則. 会話によるニュース記事伝達のための談話構造解析. 言語処理学会 第 27 回年次大会 発表論文集, 2021.
- [35] 高津弘明, 安藤涼太, 本田裕, 松山洋一, 小林哲則. 談話構造制約付きパーソナライズド抽出型要約. 言語処理学会 第 27 回年次大会 発表論文集, 2021.