

# 科学技術論文を対象とした根拠付き生成型要約システムの構築

笠原智仁 村田栄樹 河原大輔

早稲田大学理工学術院

{tomo\_k@ruri.,eiki.1650-2951@toki.,dkw@}waseda.jp

## 概要

本研究では日本語の科学技術論文を対象として、抽出型要約と生成型要約を融合した根拠付きの生成型要約システムを構築する。具体的には、論文本文から抽出型要約を作成し、それを大規模言語モデルによって言い換えることによって生成型要約を生成する。抽出型要約によって本文から抽出された文が生成された要約の根拠となり、言い換えによって自然な文章を生成することが可能となる。評価の結果、本文から直接要約を生成する生成型要約モデルと比較して、幻覚の生成が減少することを確認した。

## 1 はじめに

自動要約の手法は、原文から重要とされる文を取り出す抽出型要約と、原文を入力して要約を生成させる生成型要約の2種類に大きく分けられる。抽出型要約は誤った情報が含まれないという長所があるが、抽出された要約が文章として不自然であるという欠点がある。一方で生成型要約においては、自然な文章を生成することが可能になるという長所があるが、生成された要約が原文の内容に忠実か保証されないという欠点が存在する。このような原文に基づいていない、捏造された文章が生成される事象は幻覚 (Hallucination) と呼ばれる。

本研究ではこれらの欠点を克服するため、図1に示すような、2つの手法を組み合わせた根拠付き生成型要約システムを構築する。具体的には、原文から抽出型要約を作成し、それを大規模言語モデルによって言い換えることによって要約を生成するという2つのステップを踏む。本手法では、抽出型要約によって原文から抽出された文が生成された要約の根拠となり、言い換えによって自然な文章を生成することが可能となる。

日本語の科学技術論文を対象として実験を行い、根拠となる文を示した上で生成型要約モデルと同

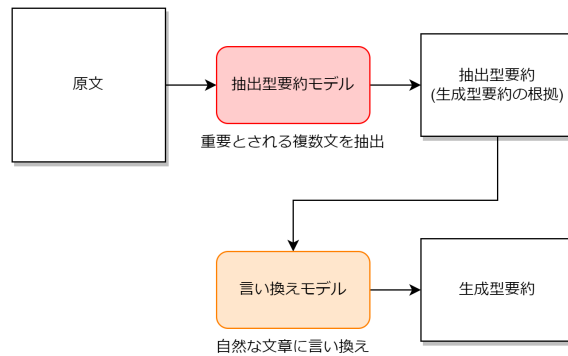


図1 根拠付き生成型要約システムのアーキテクチャ

等の精度を出すことが可能であることを示した。また、GPT-4を用いて幻覚についての評価を行った結果、生成型要約モデルと比較して幻覚が発生しづらくなるということも明らかになった。

## 2 関連研究

抽出型要約と要約型要約を組み合わせた手法の先行研究としては、Transformer [1]の考案以前からRNNを利用した手法 [2, 3]が提案されている。Transformerの考案により要約モデルの性能は向上したが、依然として抽出型と生成型それぞれの欠点は残るため、2種類を融合した手法の研究は続けられている。

抽出型から生成型への情報の受け渡しには様々な方法があり、原文から文単位で抽出した文章を生成型要約モデルの入力とする手法 [4]や、それ以外のキーワードや原文そのものも入力に含める手法 [5, 6]、離散的な抽出ではなく原文を文単位で重みづけする手法 [7]などがある。本研究では抽出型要約が生成された要約の根拠となることに着目しているため、抽出型要約は文単位の離散的な抽出であり、生成型要約への入力には抽出された原文以外は用いないという条件を満たさなければならない。

### 3 根拠付き生成型要約システム

本研究では抽出型要約と生成型要約を組み合わせることにより、根拠付き生成型要約システムを構築する(図 1)。要約の流れとしては、まず原文を抽出型要約モデルに入力し、重要とされる文を複数文抽出する。ここで抽出された文は最終的に生成される要約の根拠となる。次に、抽出された文を言い換えモデルに入力し、要約を生成する。

#### 3.1 抽出型要約モデル

抽出型要約モデルは Transformer の encoder ベースの事前学習済みモデルを、BERTSUM [8] の手法を参考に fine-tuning することで構築する。具体的には、原文の各文の先頭に CLS トークンを挿入し入力とする。次に、その埋め込みを基に抽出すべきか否かを示すスコアを出力する。

#### 3.2 言い換えモデル

言い換えモデルは Transformer の decoder ベースの事前学習済みモデルを fine-tuning することで構築する。原文から抽出した複数文をモデルに入力し、自然な要約が生成されるように学習を行う。

## 4 データセット

科学技術振興機構が本研究のために国内発行の科学技術文献より集めたデータセットを利用する。このデータセットには、本文と抄録を持つ論文(全て学会予稿)が 22,209 件含まれている。本文は 84% の文献が日本語で執筆されており、残りは英語からの機械翻訳である。一方、抄録は 5% のみが日本語であり、95% が英語からの機械翻訳である。これは、予稿がすでに抄録的な意味合いを持つため、日本語の抄録は略し、英語読者のための英語抄録のみが存在することが多いためと考えられる。

構築する根拠付き生成型要約システムは抽出型要約モデルと言い換えモデルからなる。それぞれ独立に学習を行うため、その学習データについて 4.1、4.2 節で述べる。なお、いずれのモデルの学習においても、train:valid:test=8:1:1 に分割し、test は共通のデータとなるようにした。

#### 4.1 抽出型要約モデルの学習データ

抽出型要約モデルの学習に用いる正解データは、本文の各文について抽出すべきかの 2 値分類ラベ

ルを付与することで作成する。具体的には抄録の各文に対して、本文の各文もしくは連続する 2 文との間で BLEU スコアを計算し、最もスコアが高い 1 文または連続する 2 文に「抽出」ラベルを付与する。ただし、既に抽出済みの場合には、次にスコアの高い文または連続する 2 文にラベルを付与する。連続する 2 文を考慮する理由は、本文では 2 文に分かれている内容が抄録では 1 文にまとめられるケースがあるためである。以降、このラベルが付与された本文内の文集合を「BLEU 抽出」と呼ぶ。

#### 4.2 言い換えモデルの学習データ

言い換えモデルの学習時の入力には、上記の BLEU 抽出と、学習を行った抽出型要約モデルによって抽出した文集合の、2 通りを実験して比較する。なお、後者については学習データと推論データの重複を避けるために 5-fold cross-validation によって作成する。

#### 4.3 幻覚のテストデータ

本研究では、生成した要約における幻覚の評価は GPT-4 を用いて行う。上記データセットは権利の都合上 GPT-4 へのデータ投入ができないため、投入が可能な 23 種類の資料から 10 論文ずつをランダムに選択して幻覚のテストデータを構築した。なお、この論文データには本文のみが含まれ、抄録は付属していない。

## 5 実験

抽出型要約モデルと言い換えモデルのそれぞれを独立して学習と評価を行った後に、End-to-End での評価を行う。

### 5.1 実験設定

モデル学習時のハイパーパラメータについては付録 A に記載する。

#### 5.1.1 根拠付き生成型要約システムの設定

抽出型要約モデルと言い換えモデルには、それぞれ日本語で事前学習された BigBird<sup>1)</sup> [9] と GPT-2 [10] ベースのモデル<sup>2)</sup>を用いる。

抽出型要約モデルの推論時には閾値を設定し、そ

1) <https://huggingface.co/nlp-waseda/bigbird-base-japanese>

2) <https://huggingface.co/line-corporation/japanese-large-lm-1.7b>

**表 1** 生成された抄録の評価。言い換えモデル単体、End-to-End、生成ベースラインの結果をそれぞれ示す。

設定	学習データ	推論データ	ROUGE-1	ROUGE-2	ROUGE-L	平均文数	# 繰り返し
言い換えモデル	BLEU 抽出	BLEU 抽出	0.453	0.178	0.306	4.17	87
	BigBird 抽出	BLEU 抽出	0.392	0.149	0.266	3.08	355
End-to-End	BLEU 抽出	BigBird 抽出	0.389	0.141	0.243	5.42	88
	BigBird 抽出	BigBird 抽出	0.376	0.131	0.245	3.50	294
生成ベースライン	本文	本文	0.393	0.142	0.250	4.22	79

**表 2** 抽出型要約モデルの評価。Reference は BLEU 抽出。

閾値	Accuracy	Recall	Precision	F1
0.1	0.436	<b>0.919</b>	0.270	0.398
0.2	0.621	0.731	0.329	<b>0.433</b>
0.3	0.725	0.549	0.397	0.432
0.4	0.769	0.430	0.449	0.422
0.5	<b>0.778</b>	0.382	<b>0.466</b>	0.397

**表 3** 抽出型要約モデルの評価。Reference は正解抄録。

手法	閾値	ROUGE-1	ROUGE-2	ROUGE-L
BLEU	-	0.527	0.232	0.351
BigBird	0.2	0.363	0.155	0.226
	0.3	0.421	0.167	0.255

の値を超えた文のみを抽出する。また、最低抽出文数を 3 に設定し、閾値を超えた文数が 3 文未満の場合にはスコアの高い 3 文を抽出することとした。

言い換えモデルでは、入力と正解抄録の間に SEP トークンを挿入し fine-tuning する。抄録部のみについて損失を計算する。推論時は、SEP トークンまで入力し、続きを生成することで抄録とする。

### 5.1.2 ベースラインモデルの設定

生成型要約のベースラインとして、本文から直接、抄録を生成するモデルを構築する。ベースモデルは 5.1.1 節と同じ GPT-2 を用いる。本文を入力、抄録を正解として fine-tuning する。本文、“要約:”、抄録の 3 つの文字列を連結し、抄録部のみで損失計算を行う。推論時は本文と“要約:”を入力し、続きを生成することで抄録とする。

### 5.1.3 評価指標

抽出型要約モデルの評価として、Accuracy、Recall、Precision、F1 による評価と、表層一致ベースの評価を行う。生成抄録の評価には、表層一致ベースの評価と、幻覚についての評価を行う。表層一致ベースの評価には、正解抄録との ROUGE- $\{1, 2, L\}$  を計算する。

幻覚の評価については GPT-4 に本文と抄録を入力し、幻覚を含むかどうかを判定する。OpenAI API を通じて gpt-4-turbo の Chain-of-Thought プロンプティング [11] により、幻覚があれば 1 を、なければ 0 を

**表 4** 抽出型要約モデルによる抽出と、本文と正解抄録の文数。

本文	抄録	BLEU 抽出	BigBird による抽出	
			0.2	0.3
25.37	4.14	4.45	9.88	5.72

出力させる。詳細は付録 B に示す。GPT-4 の出力を当該要約 1 文あたりに正規化し、これを評価セットで平均した値を幻覚率とする。ただし、幻覚の評価は End-to-End の設定でのみ行う。

生成では繰り返しが起こることもあるため、上記の評価は繰り返しを除去した後に行う。使用するモデルのトークンベースで、同一の 5-gram が 15 回以上出現した場合に繰り返しとみなす。

## 5.2 実験結果

### 5.2.1 抽出型要約モデルの評価

5.1.1 節の設定で fine-tuning を行った抽出型要約モデルによる抽出結果と、BLEU 抽出との間における、Accuracy、Recall、Precision、F1 による評価結果を表 2 に示す。F1 において、閾値が 0.2 と 0.3 の場合にスコアが高くなったため、以降ではこの 2 種類の閾値の結果を載せる。

抽出型要約モデルによる抽出結果と正解抄録との間における、ROUGE による評価結果を表 3 に示す。また、表 4 に各項目ごとの文数を示す。正解抄録や閾値が 0.3 の場合と比較して、閾値が 0.2 の場合は抽出される文数が多くなってしまい、ROUGE のスコアも低いことが分かる。これらの結果から、言い換えモデルの学習や End-to-End での推論には閾値が 0.3 のデータを用いることとした。

### 5.2.2 言い換えモデルの評価

言い換えモデル単体での評価結果を表 1 の上部に示す。推論時の入力は BLEU 抽出を使用する。fine-tuning 時の入力は 4.2 節で述べた 2 つの設定、すなわち BLEU 抽出と BigBird による抽出 (BigBird 抽出) で比較する。

BigBird 抽出で fine-tuning したモデルでは、繰り返



**表 5** 幻覚の評価 (End-to-End)。提案手法横の括弧内には fine-tuning に使用した抽出方法を示す。幻覚率は生成された抄録の 1 文あたりに幻覚が含まれる割合である。

モデル	# 忠実	# 幻覚	# 繰り返し	平均文数	幻覚率 (%)
BigBird 抽出	205	25	0	6.3	1.71
提案手法 (BLEU)	53	163	14	6.2	12.22
提案手法 (BigBird)	89	117	24	3.0	18.99
生成ベースライン	84	131	15	3.3	18.32

**表 6** 提案手法 (学習:BLEU 抽出、抽出閾値:0.3) と生成型ベースラインの定性的な比較。太字で BigBird によって抽出された本文を示し、幻覚となっている部分をハイライトした。

本文	提案手法 (BLEU, 0.3)	ベースライン
<p>～前略～今回、<b>血便・便中カルプロテクチン高値の精査目的に上下部消化管内視鏡検査を施行し、若年性ポリープの診断を得た一例を経験したので報告する。</b> 8 か月前からの血便を主訴に 7 か月前に当院小児科を紹介受診した。～中略～術後 3 か月後の便中カルプロテクチンは <math>14.9 \mu\text{g/g}</math> と改善を認め、血便の症状なく経過している。～中略～【結語】<b>血便・便中カルプロテクチン高値に対して炎症性腸疾患を鑑別に上下部消化管内視鏡検査を施行したが、若年性ポリープを認めEMRをおこなった。若年性ポリープの炎症細胞浸潤の病態が便中カルプロテクチン上昇に反映されることから、若年性ポリープも鑑別として挙げEMRの準備も行うことが重要である。</b></p>	<p>血便と便中のカルプロテクチンの高値を精査目的に受診した一例を経験した。内視鏡検査で若年性ポリープを認めたため、切除術を行った。若年性ポリープの炎症性腸疾患の鑑別に、血便と便中のカルプロテクチンの高値も考慮する必要がある。</p>	<p>鮮血混じりの茶褐色便をほぼ毎日認め、貧血の進行なく、腹部症状や発熱のない<b>8 か月の男児</b>の精査目的に、全身麻酔下上下部消化管内視鏡検査を施行した。便中カルプロテクチンは <math>2300 \mu\text{g/g}</math> で、炎症性腸疾患の精査目的に全身麻酔下上下部消化管内視鏡検査を施行した。上部消化管内視鏡検査では明らかな粘膜病変を認めず、下部消化管内視鏡検査では横行結腸に <math>2\text{cm}</math> 大の有茎性ポリープを認めた。EMR を施行した。若年性ポリープの診断であった。術後 3 か月の便中カルプロテクチンは <math>14.9 \mu\text{g/g}</math> と改善を認めた。</p>

しが多く ROUGE スコアも低くなっている。BigBird 抽出には、抄録として必要ではない文が含まれる可能性がある。これらがノイズとなり、学習が難しくなると考えられる。

### 5.2.3 End-to-End での評価

BigBird 抽出を言い換えモデルの推論時の入力とした End-to-End での要約を評価する。

**表層一致ベースの評価** まず、表層一致ベースの結果を表 1 の下部に示す。5.1.2 節で構築した生成型要約のベースラインの結果も含める。

5.2.2 節と同様に、BigBird 抽出を言い換えモデルの fine-tuning に使用したモデルは、繰り返しが多い。Fine-tuning 時の入力と推論時の入力の性質が似ている方がよい結果を得られると考えたが、高品質なデータでの fine-tuning がよりよい結果を得た。

また ROUGE の指標では、生成ベースラインが提案手法を若干上回った。ただし、BLEU 抽出を言い換えモデルの入力としたモデルと比較すると生成ベースラインが劣るため、ボトルネックは抽出型要約にあると考えられる。

**幻覚の評価** 表 5 に GPT-4 による幻覚の評価結果を示す。BigBird 抽出、提案手法 (学習データの異なる

2 設定)、生成ベースラインの結果をそれぞれ示す。ただし、4.3 節で言及したデータの差には留意されたい。

モデルによって出力する文数が異なるため、文あたりの幻覚率に着目すると、BLEU 抽出で fine-tuning した提案手法が生成ベースラインよりもよい結果を得た。これは根拠付き生成型要約であることによる利点である。一方、BigBird 抽出で fine-tuning した提案手法はやや生成ベースラインより劣った。これは言い換えモデルの fine-tuning 時に入力と正解抄録で参照する情報が異なるケースが多いためと考えられる。

表 6 に提案手法と生成ベースラインが生成した抄録の例を示す。ベースラインでは本文に基づかない生成が見られるが、提案手法ではより忠実な生成が多かった。

## 6 おわりに

本研究では日本語の科学技術論文を対象として、抽出型要約と生成型要約を融合した根拠付きの生成型要約システムを構築した。実験の結果、生成型要約モデルと比較して幻覚の発生確率が下がることが明らかになった。

## 謝辞

本研究は国立研究開発法人科学技術振興機構「科学技術文献の活用業務に係る自然言語処理研究および技術実証事業」の委託研究において行った。本研究に的確な助言を頂いた国立研究開発法人科学技術振興機構の鈴木慶二氏、菊井玄一郎氏に感謝する。

## 参考文献

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. **Advances in neural information processing systems**, Vol. 30, , 2017.
- [2] Yen-Chun Chen and Mohit Bansal. Fast abstractive summarization with reinforce-selected sentence rewriting. In Iryna Gurevych and Yusuke Miyao, editors, **Proceedings of the 56th ACL**, pp. 675–686, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [3] Wan-Ting Hsu, Chieh-Kai Lin, Ming-Ying Lee, Kerui Min, Jing Tang, and Min Sun. A unified model for extractive and abstractive summarization using inconsistency loss. In Iryna Gurevych and Yusuke Miyao, editors, **Proceedings of the 56th ACL**, pp. 132–141, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [4] Haoran Li, Arash Einolghozati, Srinivasan Iyer, Bhargavi Paranjape, Yashar Mehdad, Sonal Gupta, and Marjan Ghazvininejad. EASE: Extractive-abstractive summarization end-to-end using the information bottleneck principle. In Giuseppe Carenini, Jackie Chi Kit Cheung, Yue Dong, Fei Liu, and Lu Wang, editors, **Proceedings of the Third Workshop on New Frontiers in Summarization**, pp. 85–95, Online and in Dominican Republic, November 2021. Association for Computational Linguistics.
- [5] Jonathan Pilault, Raymond Li, Sandeep Subramanian, and Chris Pal. On extractive and abstractive neural document summarization with transformer language models. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, **Proceedings of the 2020 Conference on EMNLP**, pp. 9308–9319, Online, November 2020. Association for Computational Linguistics.
- [6] Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao Jiang, and Graham Neubig. GSum: A general framework for guided neural abstractive summarization. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, **Proceedings of the 2021 Conference of the NAACL-HLT**, pp. 4830–4842, Online, June 2021. Association for Computational Linguistics.
- [7] Fei Wang, Kaiqiang Song, Hongming Zhang, Lifeng Jin, Sangwoo Cho, Wenlin Yao, Xiaoyang Wang, Muhao Chen, and Dong Yu. Saliency allocation as guidance for abstractive summarization. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, **Proceedings of the 2022 Conference on EMNLP**, pp. 6094–6106, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [8] Yang Liu and Mirella Lapata. Text summarization with pretrained encoders. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, **Proceedings of the 2019 Conference on EMNLP-IJCNLP**, pp. 3730–3740, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [9] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. **Advances in neural information processing systems**, Vol. 33, pp. 17283–17297, 2020.
- [10] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. **OpenAI blog**, Vol. 1, No. 8, p. 9, 2019.
- [11] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, **Advances in Neural Information Processing Systems**, Vol. 35, pp. 24824–24837. Curran Associates, Inc., 2022.
- [12] Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansky, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. TRUE: Re-evaluating factual consistency evaluation. In Song Feng, Hui Wan, Caixia Yuan, and Han Yu, editors, **Proceedings of the Second DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering**, pp. 161–175, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [13] Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. SummEval: Re-evaluating Summarization Evaluation. **Transactions of the Association for Computational Linguistics**, Vol. 9, pp. 391–409, 04 2021.

## A ハイパーパラメータ

モデルの学習時に設定したハイパーパラメータを表 7 に示す。なお、言い換えモデルの推論では貪欲法を用いた。

ハイパーパラメータ	抽出型要約モデル	言い換えモデル
Learning Rate	1e-5	2e-5
Epoch Num	10	5
Batch Size	1	1

## B 幻覚評価の詳細

5.1.3 節で導入した、GPT-4 による幻覚評価の詳細を述べる。

### B.1 幻覚評価の手法

GPT-4 を使用した zero-shot Chain-of-Thought により、生成された抄録に幻覚が含まれるかを判定する。使用したチェックポイントは“gpt-4-1106-preview”である。

まず図 2(a) のように本文と抄録を提示し、幻覚が含まれるか記述する。“source”と“target”にそれぞれ本文と生成された抄録を代入する。次にプロンプト 1 とそれに対する GPT-4 の出力に加えて、図 2(b) のプロンプト 2 を入力することで 2 値の評価を得る。

[要約] に含まれる情報が [本文] と矛盾しないか判断してください。もし [要約] が [本文] に基づいていない箇所を含む場合、その [要約] の部分と判断材料となった [本文] の該当箇所を示して理由を説明してください。\*\*\*ただし、[本文] に含まれるすべての事実が [要約] で言及される必要はありません。\*\*\*

[本文]  
{source}

[要約]  
{target}

(a) プロンプト 1

以上の情報を踏まえて、[要約] に含まれる情報が [本文] と矛盾しているか判断してください。矛盾している場合は 1 を、矛盾していない場合は 0 を出力してください。数字のみを出力してください。理由などは必要ありません。

(b) プロンプト 2

図 2 幻覚評価に使用したプロンプト

### B.2 幻覚評価の性能

この評価の性能を検証するため、TRUE [12] で 2 値化された SummEval [13] でテストを行う。新聞記事とモデルの出力した要約に対して、幻覚が含まれていないか人手でアノテーションされたデータである。このデータに対して B.1 節の手法で評価した結果を表 8 に示す。

表 8 ラベル付き要約データでの幻覚評価器の評価結果

Accuracy	Recall	Precision	F1
0.720	0.790	0.693	0.738