

広告文生成における多様性と広告品質の関係

青木洋一^{1,2} 村上聡一郎³ 本多右京³ 加藤明彦³

¹ 東北大学 ² 理化学研究所 ³ 株式会社サイバーエージェント

youichi.aoki.p2@dc.tohoku.ac.jp,

{murakami_soichiro,honda_ukyo,kato_akihiro}@cyberagent.co.jp

概要

広告文生成において、幅広いユーザーの関心を引き、同じ広告に繰り返し触れることで引き起こされる広告疲れを防ぐためには、多様な広告文を生成することが必要不可欠である。しかしながら、大規模言語モデルを用いて多様な広告文を生成する場合、広告文の多様性を向上させることが広告品質にどのような影響を与えるのかはよく分かっていない。要約や機械翻訳などのタスクではこうした多様性と品質の関係が検証されてきたが、広告文生成はこれらタスクと比較して文体や品質の評価軸が大きく異なり、多様性と品質の関係は非自明である。そこで本研究では、デコーディング手法、ハイパーパラメータ、入出力形式、モデルの数といったいくつかの要因を変化させることで、日本語広告文生成における多様性と広告品質の関係を明らかにした。

1 はじめに

広告は企業が自社の製品やサービスを幅広い層に訴求するために必要不可欠な手段である。自然言語処理分野ではこうした広告の作成需要に応えるため、言語モデルを活用した広告文生成に関する研究が進められている [1, 2, 3, 4]。

広告文生成において多様性は重要な指標である。同一の広告を繰り返し提示すると、ユーザが飽きや広告疲れを引き起こす可能性がある [5, 6]。また、多様な広告文を提示することで、より幅広い顧客層に訴求できる可能性が高まる [7]。したがって、多様な広告文を自動生成することが強く求められている。

しかしながら、大規模言語モデル (LLM) を用いて多様な広告文を生成する場合、beam search [8, 9, 10] や sampling [11, 12, 13] といったデコーディング手法による多様性の増加が、広告品質などの他評価指標にどのような影響を及ぼすかは十分に分かってい

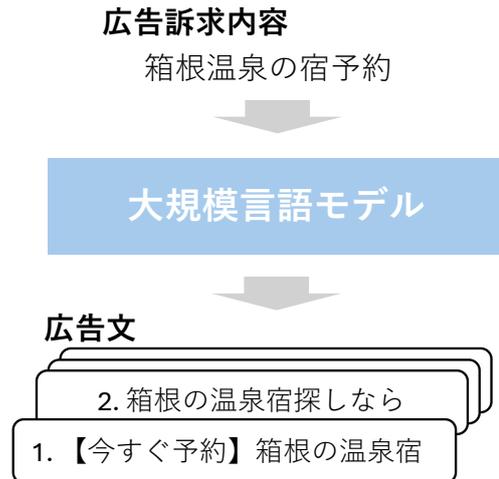


図1 大規模言語モデルに対する入出力例。入力は広告訴求内容であり、出力は複数の広告文である。

ない。先行研究によれば、多様性と品質の関係はタスク毎に大きく異なることが報告されている [14]。例えば、要約や機械翻訳では多様性と品質の間にトレードオフの関係が見られる一方で、物語生成では両者は相互依存の関係であることが観察されている。広告文生成はこうしたタスクと以下の点で異なる。

1. 文法的な誤りがある程度許容される一方で、広告効果や文の長さといった要因がより重視される [15]。
2. スローガンやバナー広告などでは限られた文字数で効果的にメッセージを伝える必要があり、記号やキーワードを活用した独自の表現が多く見られる [16]。

こうした評価基準や文体の違いから、広告文生成における多様性と品質の関係は非自明であり、体系的な検証が必要である。本研究では、デコーディング手法、ハイパーパラメータ、入出力形式、モデルといった複数の要因を変化させることで、日本語広告文生成における多様性と広告品質の関係を調査

した。

2 問題設定

2.1 入出力

本研究では、Murakami 等 [16] の設定に基づき、LLM に対して、日本語広告文データセットである CAMERA [17] の広告文を広告訴求内容として与え、5つの広告文を生成させる (図 1 参照)。ここでいう広告訴求内容とは、広告がユーザーに対して強調する製品・サービスの主要な価値や魅力を指す。また、文脈内学習のために LLM には 3 組の入出力例を提示する。

2.2 デコーディング手法

実験では sampling [13, 11]、beam search [8, 9, 10] および Diverse Beam Search (DBS) [18] といった典型的なデコーディング手法が多様性に与える影響を調査する。また、翻訳やキャプション生成などのタスクで多様性と品質のトレードオフを改善することが報告されている Diverse Minimum Bayes Risk (DMBR) および k-medoids Minimum Bayes Risk (KMBR) [19] などの MBR デコーディング手法が、広告文生成で多様性に与える影響も調査する。

2.3 評価

本研究では、生成された 5 つの広告文間の多様性と、各広告文の品質を測定する。

2.3.1 広告文の多様性

広告文の多様性は意味的多様性と表層的多様性の 2 つに分類できる [20]。意味的多様性は、広告文の訴求内容の違いに関する観点であり、表層的多様性は訴求内容が同一の文間での表層的な文字列の違いに関する観点である。

一般に広告制作は次の 2 つの段階に大別される。

- (1) 効果的な訴求内容を探索する段階：広告配信初期における訴求の検証段階
- (2) 特定の訴求内容に対して効果的な表現を探索する段階：有効な訴求が判明した後の表現最適化段階

(1) では、競合他社の訴求内容や商品の特徴を踏まえて複数の訴求案を設定し、それぞれに基づく広告文を生成する。(2) では、特定の訴求を固定したう

えで表層的に多様な表現を生成する。したがって、いずれの段階においても広告文は指定された訴求内容に沿って作成される必要がある。このような訴求制約下での多様性を実現するために、本研究では Murakami 等 [16] に従い、表層的多様性に焦点を当てる。表層的に多様な広告文の具体例は以下のとおりである。

- サービス：携帯食
- 訴求内容：短時間で健康的な食事
- 広告文例：
 - 忙しく働く人のための即効エネルギー
 - 時間を節約、健康はそのまま
 - 早く健康的、これ一食で

表層的多様性の評価には、出力文同士の n -gram 一致に基づくペアワイズ BLEU 値 [21] を用いる。具体的には、多様性指標として $1 - \text{BLEU}$ を測定する。この値は 0 から 1 の範囲を取り、値が大きいほど多様性が高いことを示す。

2.3.2 広告品質

広告品質は、先行研究で定義された代表的な 3 つの指標——広告効果、一貫性、および受容性——に基づいて評価する [15]。1 つの入力に対して生成された 5 文の平均値を広告品質と定義する。

広告効果 先行研究 [1, 22, 17] に従い、過去の配信履歴に基づくクリック率 (CTR) をもとに、顧客行動をシミュレーションし広告効果を測定する。¹⁾ 本研究では、CAMERA データセットに含まれる人手で作成された広告文の予測 CTR と、LLM によって生成された広告文の予測 CTR 比 (生成文 CTR / 参照文 CTR) を評価指標として使用する。

一貫性 入力された広告内容と異なる訴求内容を生成すると、広告主に不利益をもたらす可能性がある。例えば、「有料」を「無料」と誤って生成した場合、虚偽広告とみなされるおそれがある。そのため、入力訴求内容と生成広告文の訴求内容の一貫性を、BERTScore [23] を用いて評価する。

受容性 広告プラットフォームでは文字数制限が設けられている場合が多い。本研究では、全角 15 文字または半角 30 文字以内に収まる場合を広告文として「受容可能」と定義する。²⁾

1) CTR 予測モデルである極予測 TD を使用した。本モデルの予測値は実際の CTR と整合している。cf. <https://cyberagent.ai/products/>

2) Google 広告など日本語広告プラットフォームの一般的な制約に基づく。cf. <https://support.google.com/>

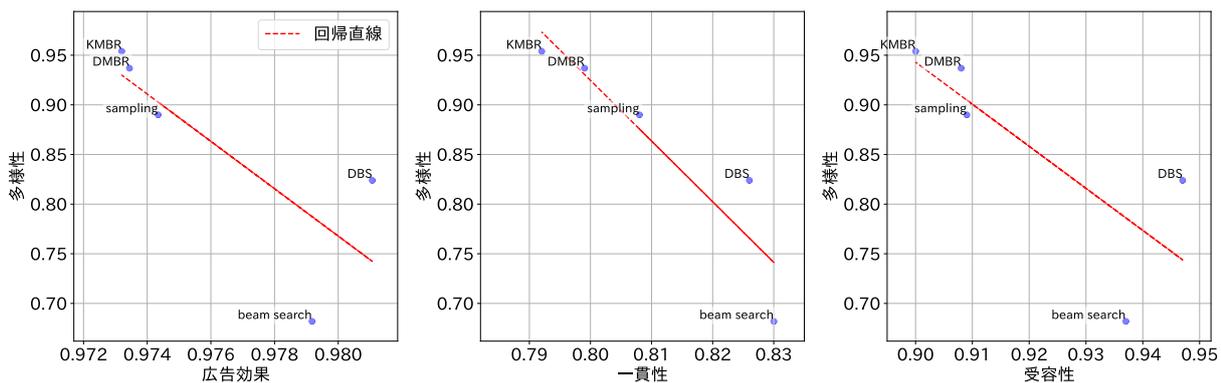


図2 広告文生成における多様性と広告品質の関係。図中の各点は、LLMに任意のデコーディング手法を適用して広告文を生成させた際の生成広告文の多様性と広告品質の値を表す。

3 実験

本研究では、calm3-22b-chat (calm3) [24]、Llama-3-ELYZA-JP-8B (ELYZA) [25]、Mistral-Small-24B-Instruct-2501 (Mistral) [26]、Llama-3.1-Swallow-70B-Instruct-v0.3 (Swallow) [27, 28, 29]、および GPT-4o [30] を使用した。複数のモデル間で類似した傾向が確認されたため、本論文では calm3 の結果のみを本文中に示し、その他のモデルに関する結果は付録 B に記載した。

3.1 異なるデコーディング手法間における多様性と広告品質の関係

図2に、異なるデコーディング手法を用いた時の生成広告文の多様性と品質をプロットした。図2に示すように、多様性と広告品質の間にトレードオフの関係が見られた。この傾向は、対話生成や物語生成といった他の生成タスクで報告されている関係 [14] とは異なるものである。これは、対話や物語生成はより創造的かつ動的な性質を持つものに対し、広告文生成は目標志向型のタスクであり、明確な訴求内容や制約条件の下で生成が行われるからだと考えられる。一方で、beam search は DBS と比較して、広告効果や受容性を維持しつつ多様性を向上させる結果が観察された。また、例外として、GPT-4o では多様性と受容性の双方を改善するバランスの取れた結果が得られた。これらの傾向は人手評価の結果とも整合していた (付録 C 参照)。

3.2 デコーディング手法のパラメータ変更と多様性・広告品質への影響

sampling および beam search といったデコーディング手法のパラメータを変更した時の多様性と広告

品質の関係を表1に示した。KMBR は代表的な変更可能なパラメータが代表候補数 K であり、出力文数を固定している本実験設定では変化が少ないと思われることから、本実験の対象外とした。

sampling と DMBR に関しては、確率閾値 P や温度パラメータ T 、多様性ペナルティ D を大きくするほど多様性が向上し、品質が低下する傾向が確認できた。一方、beam search では、ビーム幅 W を10まで拡大する場合は、一貫性を保ちながら多様性が向上した。これは、機械翻訳に関する先行研究 [31] とは対照的である。この違いは、広告文生成が機械翻訳よりも短いテキストを生成するタスクであるためと考えられる。また、DBS では、ビーム幅を増加させると、多様性、広告効果、受容性のいずれも低下する結果が得られた。

3.3 Few-shot 事例数の変化が多様性および広告品質に与える影響

表2に、few-shot 学習における few-shot 事例数の違いが多様性および広告品質に与える影響を示した。KMBR と DMBR は sampling 後の生成広告文を評価して最終的な生成文を決定する手法であり、few-shot 事例数を変化させた時の多様性と品質の関係は sampling と同様の傾向になると考えられるため、本実験からは除外した。

sampling においては、few-shot 事例数を増やすことで広告品質は向上するが、多様性が減少する傾向が見られた。一方、beam search と DBS では、few-shot 事例数を増やすことで多様性と広告品質の双方が向上した。これは、LLM に few-shot 事例として多様な例を与えることで、出力確率の高い語彙が多様化したためと考えられる。

表 1 各デコーディング手法のパラメータを変更した際の多様性および広告品質。sampling では確率閾値 (P) および温度パラメータ (T) を変更し、beam search ではビーム幅 (W) を変更した。DBS では、ビーム幅 (W) およびグループ数 (G) の両方を変更した。DMBR では多様性ペナルティ (D) を変更した。

sampling				
パラメータ	多様性	広告効果	一貫性	受容性
$P=0.5, T=1$	0.515	0.981	0.832	0.940
$P=1, T=1$	0.890	0.974	0.808	0.909
$P=1, T=1.5$	0.967	0.967	0.773	0.809
$P=1, T=2$	0.991	0.963	0.740	0.619
beam search				
パラメータ	多様性	広告効果	一貫性	受容性
$W = 5$	0.682	0.979	0.803	0.937
$W = 7$	0.697	0.980	0.832	0.939
$W = 10$	0.702	0.980	0.831	0.932
$W = 12$	0.699	0.980	0.832	0.933
DBS				
パラメータ	多様性	広告効果	一貫性	受容性
$W, G = 5$	0.824	0.981	0.826	0.947
$W, G = 7$	0.814	0.980	0.825	0.941
$W, G = 10$	0.762	0.980	0.826	0.927
$W, G = 12$	0.739	0.980	0.827	0.913
DMBR				
パラメータ	多様性	広告効果	一貫性	受容性
$D = 0.1$	0.767	0.978	0.826	0.922
$D = 0.5$	0.927	0.974	0.802	0.904
$D = 1$	0.954	0.973	0.792	0.9
$D = 2$	0.957	0.973	0.79	0.897

3.4 複数モデルを組み合わせた広告文生成

これまでの実験は単一のモデルに複数文を出力させる設定で行われた。そこで新たに、異なるモデルがそれぞれ独立に広告文を生成する場合の多様性と広告品質の関係を調査した。具体的には、5つの異なるモデルそれぞれが1文ずつ独立に広告文を生成し、生成された5つの広告文間の多様性と広告品質を測定した。この生成手法を以降では five-model とする。§ 2.1 で述べた単一モデルによる5文生成と、five-models を比較した結果を表 3 に示した。five-models の広告品質は各モデルの平均性能とほぼ同等であったが、多様性は他の手法より高い結果となった。これは、複数のモデルを組み合わせたことが単純により多様な広告文を生成できる可能性を示

表 2 Few-shot 事例数変更時の多様性および広告品質

sampling				
事例数	多様性	広告効果	一貫性	受容性
3	0.890	0.974	0.808	0.909
9	0.872	0.985	0.822	0.975
15	0.876	0.987	0.825	0.982
beam search				
事例数	多様性	広告効果	一貫性	受容性
3	0.682	0.979	0.830	0.937
9	0.738	0.989	0.845	0.982
15	0.747	0.992	0.849	0.987
DBS				
事例数	多様性	広告効果	一貫性	受容性
3	0.824	0.981	0.826	0.947
9	0.832	0.990	0.838	0.987
15	0.830	0.991	0.842	0.991

唆している。

表 3 モデル毎の多様性および広告品質

モデル	多様性	広告効果	一貫性	受容性
calm3	0.890	0.974	0.808	0.909
ELYZA	0.877	0.970	0.781	0.864
Mistral	0.886	0.984	0.833	0.877
Swallow	0.868	0.962	0.795	0.503
GPT-4o	0.777	0.990	0.823	0.992
5 models	0.929	0.975	0.808	0.827

4 おわりに

本研究では、日本語広告文生成における多様性と品質の関係を調査した。その結果、異なるデコーディング手法間で多様性と広告品質の間にトレードオフが確認できた。また、sampling では、確率閾値 P や温度パラメータ T 、few-shot 事例数を変化させた際に多様性と品質の間にトレードオフ関係が生じた一方で、beam search ではビーム幅 W や few-shot 事例数を増やした場合、多様性と品質の双方が向上する結果が得られた。これらの知見は、広告文生成においてデコーディング手法を慎重に検討する必要があることを示している。さらに、異なるモデルの出力を組み合わせることで、広告品質を損なわずに多様性を向上させる有望な方向性も示された。本研究の成果は、広告文生成における多様性と品質の両立に関する今後の研究に大きく貢献すると考えられる。

謝辞

本研究は、JST 次世代研究者挑戦的研究プログラム JPMJSP2114 の支援を受けたものです。

参考文献

- [1] J. Weston Hughes, Keng-hao Chang, and Ruofei Zhang. Generating better search engine text advertisements with deep reinforcement learning. In **KDD 2019**, pp. 2269–2277.
- [2] Hidetaka Kamigaito, Peinan Zhang, Hiroya Takamura, and Manabu Okumura. An empirical study of generating texts for search engine advertising. In **NAACL-HLT 2021**, pp. 255–262.
- [3] Konstantin Golobokov, Junyi Chai, Victor Ye Dong, Mandy Gu, Bingyu Chi, Jie Cao, Yulan Yan, and Yi Liu. Deepgen: Diverse search ad generation and real-time customization. In **EMNLP 2022**, pp. 191–199.
- [4] Soichiro Murakami, Sho Hoshino, and Peinan Zhang. Natural language generation for advertising: A survey. **CoRR**, Vol. abs/2306.12719, , 2023.
- [5] Corneilia Pechman and David W. Stewart. Advertising Repetition: A Critical Review of Wearin and Wearout. **Current issues and research in advertising**, Vol. 11, No. 1-2, pp. 285–329, 1988.
- [6] Susanne Schmidt and Martin Eisend. Advertising Repetition: A Meta-Analysis on Effective Frequency in Advertising. **Journal of Advertising**, Vol. 44, No. 4, pp. 415–428, 2015.
- [7] Penghui Wei. Ideation: Diversifying ad text generation in sponsored search with latent diffusion. In **WWW 2025**, pp. 1394–1397, 2025.
- [8] Bruce T. Lowerre. **The harpy speech recognition system**. PhD thesis, Carnegie Mellon University, 1976.
- [9] Alex Graves. Sequence transduction with recurrent neural networks. **CoRR**, Vol. abs/1211.3711, , 2012.
- [10] Alexander M. Rush, Yin-Wen Chang, and Michael Collins. Optimal beam search for machine translation. In **EMNLP 2013**, pp. 210–221.
- [11] David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. A learning algorithm for boltzmann machines. **Cogn. Sci.**, Vol. 9, No. 1, pp. 147–169, 1985.
- [12] Angela Fan, Mike Lewis, and Yann N. Dauphin. Hierarchical neural story generation. In **ACL 2018**, pp. 889–898.
- [13] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In **ICLR 2020**.
- [14] Clara Meister, Gian Wiher, and Ryan Cotterell. On decoding strategies for neural text generators. **Trans. Assoc. Comput. Linguistics**, Vol. 10, pp. 997–1012, 2022.
- [15] Peinan Zhang, Yusuke Sakai, Masato Mita, Hiroki Ouchi, and Taro Watanabe. Adtec: A unified benchmark for evaluating text quality in search engine advertising. **CoRR**, Vol. abs/2408.05906, , 2024.
- [16] Soichiro Murakami, Peinan Zhang, Hidetaka Kamigaito, Hiroya Takamura, and Manabu Okumura. AdParaphrase: Paraphrase dataset for analyzing linguistic features toward generating attractive ad texts. In **Findings of the NAACL 2025**, pp. 1426–1439.
- [17] Masato Mita, Soichiro Murakami, Akihiko Kato, and Peinan Zhang. Striking gold in advertising: Standardization and exploration of ad text generation. In **ACL 2024**, pp. 955–972.
- [18] Jiwei Li, Will Monroe, and Dan Jurafsky. A simple, fast diverse decoding algorithm for neural generation. **CoRR**, Vol. abs/1611.08562, , 2016.
- [19] Yuu Jinnai, Ukyo Honda, Tetsuro Morimura, and Peinan Zhang. Generating diverse and high-quality texts by minimum bayes risk decoding. In **Findings of the ACL 2024**, pp. 8494–8525, 2024.
- [20] Guy Tevet and Jonathan Berant. Evaluating the evaluation of diversity in natural language generation. In **EACL 2021**, pp. 326–346, 2021.
- [21] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In **ACL 2002**, pp. 311–318.
- [22] Yanwu Yang and Panyu Zhai. Click-through rate prediction in online advertising: A literature review. **Inf. Process. Manag.**, Vol. 59, No. 2, p. 102853, 2022.
- [23] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with BERT. In **ICLR 2020**.
- [24] Ryosuke Ishigami. cyberagent/calm3-22b-chat, 2024.
- [25] Masato Hirakawa, Shintaro Horie, Tomoaki Nakamura, Daisuke Oba, Sam Passaglia, and Akira Sasaki. elyza/llama-3-elyza-jp-8b, 2024.
- [26] Mistral AI team. Mistral small 3. **Mistral AI News**, 2025.
- [27] Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. Continual pre-training for cross-lingual llm adaptation: Enhancing japanese language capabilities. In **Proceedings of the First Conference on Language Modeling**, COLM, 2024.
- [28] Naoaki Okazaki, Kakeru Hattori, Hirai Shota, Hiroki Iida, Masanari Ohi, Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Rio Yokota, and Sakae Mizuki. Building a large japanese web corpus for large language models. In **Proceedings of the First Conference on Language Modeling**, COLM, 2024.
- [29] Youmi Ma, Sakae Mizuki, Kazuki Fujii, Taishi Nakamura, Masanari Ohi, Hinari Shimada, Taihei Shiotani, Koshiro Saito, Koki Maeda, Kakeru Hattori, Takumi Okamoto, Shigeki Ishida, Rio Yokota, Hiroya Takamura, and Naoaki Okazaki. Building instruction-tuning datasets from human-written instructions with open-weight large language models, 2025.
- [30] OpenAI. Gpt-4o system card. **CoRR**, Vol. abs/2410.21276, , 2024.
- [31] Philipp Koehn and Rebecca Knowles. Six challenges for neural machine translation. In **NMT@ACL 2017**, pp. 28–39, 2017.
- [32] Karl Pearson and Francis Galton. Vii. note on regression and inheritance in the case of two parents. **Proceedings of the Royal Society of London**, Vol. 58, No. 347-352, pp. 240–242, 1895.

A 詳細な実験設定

本研究では、CAMERA データセットから 798 件の最大長が半角 30 文字または全角 15 文字の広告文を抽出した。

sampling では、1 回の推論につき 1 文の広告文を出力し、これを 5 回繰り返した。一方、beam search では 1 回の推論で上位 5 候補を出力とした。各手法で得られた 5 文の広告文に対して、多様性および品質を評価した。各デコーディング手法におけるパラメータ設定は表 4 に示す。LLM による広告文生成には NVIDIA A100 (80 GB) GPU を使用した。

表 4 デコーディング手法のパラメータ
デコーディング手法

デコーディング手法	パラメータ
sampling	top_p=1.0, temperature=1.0
beam search	num_beams = 5
DBS	num_beams = 5, num_beam_groups = 5, diversity_penalty=1.0
DMBR	div_pen = 1.0

B 他モデルでの結果

Llama-3-ELYZA-JP-8B (ELYZA)、Mistral-Small-24B-Instruct-2501 (Mistral)、Llama-3.1-Swallow-70B-Instruct-v0.3 (Swallow)、および OpenAI の GPT-4o (gpt-4o-2024-08-06, 以下 GPT-4o) における結果を、表 5 にそれぞれ示す。いずれのモデルにおいても、多様性と広告品質の間にトレードオフ関係が一貫して観察された。ただし、GPT-4o のみ多様性の向上と受容性の両立を達成していた。

C 人手による広告品質評価

本研究で使用した自動評価指標が人手評価とどの程度整合しているのかを評価した。具体的には、5 名のアナテータが、生成された広告文 1,000 件を人手評価した。

広告効果については、アナテータに参照文と生成広告文を与え、最も魅力的であると感じる文を選択してもらった。広告効果と魅力度は完全には同一の指標ではないが、人手による魅力度評価と自動指標による広告効果評価の一致率は 61%であった。

また、一貫性については、アナテータが生成広告文が入力文の意味を保持しているかを判定した。一

表 5 ELYZA, Mistral, Swallow, GPT-4o を用いた時の広告文生成における異なるデコーディング手法間の多様性と広告品質の関係

ELYZA				
手法名	多様性	広告効果	一貫性	受容性
sampling	0.877	0.970	0.781	0.864
beam search	0.735	0.972	0.799	0.874
DBS	0.838	0.971	0.792	0.893
DMBR	0.956	0.968	0.765	0.841
KMBR	0.935	0.973	0.799	0.906
Mistral				
手法名	多様性	広告効果	一貫性	受容性
sampling	0.886	0.984	0.833	0.877
beam search	0.557	0.991	0.878	0.943
DBS	0.780	0.988	0.862	0.904
DMBR	0.957	0.979	0.807	0.801
KMBR	0.934	0.973	0.799	0.905
Swallow				
手法名	多様性	広告効果	一貫性	受容性
sampling	0.868	0.962	0.795	0.503
beam search	0.520	0.963	0.808	0.538
DBS	0.767	0.964	0.804	0.567
DMBR	0.941	0.96	0.781	0.454
KMBR	0.934	0.955	0.764	0.400
GPT-4o				
手法名	多様性	広告効果	一貫性	受容性
sampling	0.777	0.990	0.823	0.992
DMBR	0.889	0.989	0.816	0.992
KMBR	0.875	0.990	0.821	0.993

貫性は相対評価ではないため、人手評価スコアと自動評価スコア間のピアソン相関係数 r [32] を測定した。その結果、ピアソン相関係数 r は 0.55、 p 値は 0.001 未満であり、人手評価と自動評価の間に中程度の相関が認められた。

D AI アシスタントの利用

本論文の執筆および実験用ソースコードの作成にあたり、AI アシスタント (GPT-4o や GitHub Copilot) を使用した。しかし、その利用はコード補完、文章の編集、表の作成に限定されており、内容はすべて著者らの独自のアイデアに基づくものである。