
On LLM Augmented AB Experimentation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Automated experimentation methods to evaluate user preferences and engagement
2 is a key cornerstone in the current digital landscape. Most such systems rely on
3 marketers and creators to design the content before deployment. However, with
4 the advent of Large Language Models (LLMs) the feedback cycle is considerably
5 shortened while the experimentation space expands significantly, necessitating
6 novel and efficient ways to assess user engagement. In this paper, we experiment
7 with using LLMs as simulators or treatment raters in an A/B testing application
8 without running an A/B test.

9 1 Introduction

10 Widespread adoption of mobile devices and increased internet access has led to a significant increase
11 in digital content consumption. To maximize customer engagement, businesses constantly aim to
12 optimize the content and user experience. For example, news media industries constantly strive to
13 come up with attractive headlines and cover images (Coenen, 2019) to drive customer engagement.
14 The standard practice to find attractive headlines is to use A/B testing. However, this is inefficient
15 for applications surrounding social-media, news and related sectors; as news and trends have short
16 lifetimes and might become irrelevant by the time a standard A/B test finishes. This problem is further
17 aggravated due to significant democratization of content creation, which has led to shorter feedback
18 cycles and increasing amount of content which needs to be experimented. Thus, in industries, where
19 newer content constantly comes up, there is a great need for more-efficient engagement evaluation.

20 Large language models (LLMs) have been demonstrated to have significant potential for processing
21 natural language text, following human instructions and generating high-quality responses (OpenAI,
22 2024). This has spurred their use in many applications such as tool learning (Qu et al., 2024) and
23 information retrieval (Zhu et al., 2024b). Given that LLMs have even demonstrated the ability to
24 mimic human preferences and behavior in a variety of consumer research tasks (Li et al., 2024; Brand
25 et al., 2023); a natural question is 'how useful LLMs can be for content optimization?'

26 **Contribution** In this work, we focus on the problem of using Large Language Models (LLMs) to
27 bypass current A/B testing practices. Specifically we focus on using LLMs to identify appealing
28 content. For concreteness, we will consider writing headlines for articles as our running example. As
29 such we will use the terms content/article/prompt and the terms treatment/headline interchangeably.
30 LLMs can be used in multiple ways for the purpose of rating treatments. In this paper, we explore
31 in-context learning, embedding based methods, and generative model based evaluation using two
32 benchmark datasets from real-life A/B tests. In our experiments, we find that using LLMs as few-shot
33 learners for treatment rating is significantly less effective than training models using LLM-based
34 representations. The accuracy of using in-context learning is only slightly higher than random
35 guessing. Furthermore, for methods which use LLM-based embeddings, the accuracy is not high
36 enough to be used as a standalone treatment evaluation method. Finally, we tried to use a generative

37 approach by fine-tuning the LLM to produce more engaging headlines. We found using such a tuned
38 generative model to be a more promising methodology for rating headlines.

39 2 Preliminaries

40 2.1 Learning from A/B Test

41 The language model is considered as a policy function π which observes a prompt x and produces
42 a textual response a by sampling from a distribution $y \sim \pi(\cdot | x)$. We are given a dataset of
43 $\mathcal{D}_{\text{pref}} = \{(x, a^+, a^-)\}$ of prompts and labeled response pairs. Here, a^+ is a positive response and a^-
44 is a negative response. Consider the example of A/B testing different summaries or headlines for a
45 given content. The preference data is obtained by exposing the incoming traffic to one of two possible
46 treatments/headlines (a or b) and the effective engagement (measured as clicks, screen time or any
47 other chosen metric) was monitored. The option with higher engagement is considered as the positive
48 sample a^+ while the other is considered as a^- .

49 **Offline RLHF** (Christiano et al., 2017; Ye et al., 2024a; Ouyang et al., 2022) deals with the problem
50 of aligning a policy network, using $\mathcal{D}_{\text{pref}} = \{(x, a^+, a^-)\}$. Given the context/prompt x , a pair of
51 outputs are sampled from $\pi_{\text{ref}}(\cdot | x)$ and then arranged as per preference function (typically implicitly
52 given by human annotation). RLHF methods (Christiano et al., 2017; Ouyang et al., 2022) seek to
53 obtain a policy $\hat{\pi}$ that is aligned with the preference data. This is done, by first estimating a reward
54 function r from $\mathcal{D}_{\text{pref}}$ using maximum likelihood. Then one uses RL based optimization methods like
55 PPO to maximize the learnt reward with an additional regularization term.

$$\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\pi} \left[r(x, y) - \beta \log \frac{\pi(a | x)}{\pi_{\text{ref}}(a | x)} \right].$$

56 **Overoptimization** The phenomena of overoptimization and reward-hacking in alignment literature
57 is well documented (Guo et al., 2024; Song et al., 2024). This problem can be alleviated when access
58 to the underlying system is available, as data with the policy can be collected from the policy as it
59 gets optimized (Gao et al., 2024; Guo et al., 2024). However, in the context of A/B testing, these
60 methods are inapplicable, as the only way to collect data from the newer policy is deploying it in the
61 field, i.e. another A/B test which defeats the purpose of using LLMs to bypass A/B testing.

62 2.2 Related Work

63 Researchers are increasingly trying to utilize LLMs for emulating human behaviour (Ziems et al.,
64 [n. d.]; Kim and Lee, 2023; Park et al., 2023). The idea of using AI agents to simulate users has a long
65 history of research in information systems (Carterette et al., 2011; Mostafa et al., 2003). LLM based
66 user simulators has been studied for evaluating task-oriented dialogue systems and recommender
67 systems (Balog and Zhai, 2023, 2024). Chen et al. (2024) have demonstrated the potential of using
68 self-play between LLMs for developing recommendation systems. Recent works have also suggested
69 using LLM based models to warm start bandit based methods (Ye et al., 2024b) for A/B testing.
70 However, concerns about the reliability of such simulations have also been raised (Zhu et al., 2024a).

71 3 Our Work

72 3.1 Direct Evaluation with LLM

73 LLMs have proven themselves to be good as both *embedding models* (Ethayarajh, 2019) and *task*
74 *learners* (Brown et al., 2020). We consider both of these possible ways to develop LLM based
75 baselines for rating treatments/headlines.

- 76 • **Direct Prompting:** We treat the LLM as an evaluator, provide it the article in the prompt
77 and instruct it to rate the different headlines as more engaging. This effectively uses the
78 LLM as a zero-shot classifier, and can directly measure the accuracy. We call this method
79 *PromptOnly*.

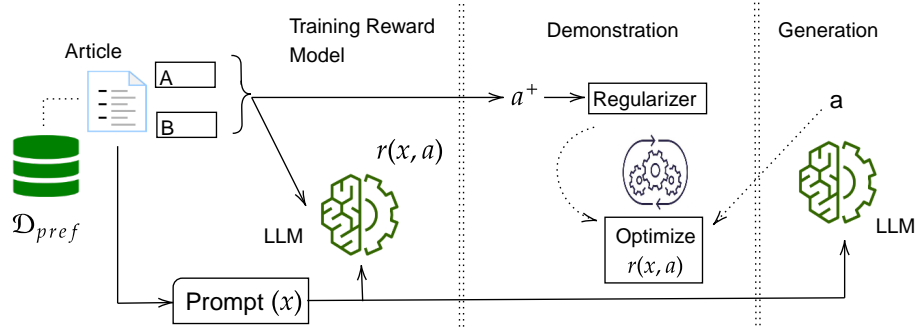


Figure 1: Overview of the proposed generative approach. The reward model r is obtained by tuning an LLM on the preference data \mathcal{D}_{pref} , which consist of tuples of contexts/articles along with two treatments arms (a^+ , a^-). Given a prompt x (which includes the context/article along with instructions) the generator LLM produces an output a . The pair x, a^+ is considered as a demonstration for the generator to match and improve using the reward model r .

- 80 • In-context learning: Similar to direct prompting except that the LLM is also provided with a
- 81 few in-context examples (or demonstrations) to learn from and choose the correct answer.
- 82 We refer to these as *ICL* methods.
- 83 • Blackbox Embedding: We train an MLP based which used the LLM embeddings of the
- 84 combined text of the article and headline to pick the better answer.
- 85 • Finetuning: We fine-tune an opensource LLM based on the data, to pick the better answer.
- 86 This is similar to the blackbox embedding approach, except that has a frozen LM, whereas
- 87 we allow the LM to be updates. Furthermore as both our compute resources and the amount
- 88 of data is limited, we take LORA (Hu et al., 2022) approach. We call these *FT* methods.

89 In the experiment section, we will present the results from all of these methods. We found that
 90 prompting and in-context learning based methods are significantly worse (< 65%) than fine-tuning
 91 based approaches (~80% accuracy). These results are qualitatively in line with other recent works
 92 focusing on using prompting and ICL based methods to classify content (Zhou et al., 2024).

93 3.2 Generative Evaluation with LLM

94 We also propose a method based on finetuning an LLM to generate engaging treatment arms (or
 95 content) using the results from the logged A/B testing data. Note that the goal here is not necessarily
 96 to use the generative model to generate new treatments, but instead use it to rate content based on
 97 model likelihood. We use RLHF (Ouyang et al., 2022) as our starting point. However we found that
 98 this model can overfit easily and is only slightly better than few-shot learning based approach. As
 99 such we modify the standard RLHF procedure to address the overfitting caused by nuances specific
 100 to the A/B testing. The overall schematic is presented in Figure 1 Most of the ideas in our approach
 101 can also be applied with other learning paradigms such as DPO (Rafailov et al., 2024), and CPO (Xu
 102 et al., 2024).

103 **Ensemble reward model** In the direct evaluation methods, we found the ensemble model from
 104 GPT embeddings to be a cheap and accurate model in predicting positive treatment arms. As
 105 such we leveraged GPT embeddings to train a reward model. However, following (Coste et al.,
 106 2024), we created an ensemble model E with different subsets of the data to help reduce reward
 107 overoptimization.

108 **Regularizing Objectives** Compared to the standard RLHF framework of (Ouyang et al., 2022) we
 109 make the following changes to the loss objective:

- 110 • We include an additional term of $\frac{\pi(a|x)}{\pi_{ref}(a|x)}$ as a regularizer in the objective. This terms
- 111 more strongly penalizes deviations of π from π_{ref} than just the KL divergence. An astute
- 112 reader might also note that this term is equivalent to regularizing with the order-2 Tsallis
- 113 divergence.

114 • We also prevent the model from exploring a space which might be far from the space where
 115 the reward model is certain. Since \mathcal{D}_{pref} forms the training data for the ensemble reward
 116 model, a y which is too far from the training data i.e. $P(y|x) \ll P(\mathcal{D}_{pref}|x)$ is low, can
 117 be considered to be a situation where the reward model is unreliable. Penalizing with the
 118 corresponding density ratio prevents the model from going too far out of the training support
 119 of the reward model. In our experiments, we estimate this ratio via prompting. Specifically,
 120 we provide GPT-3 with a prompt that describes the articles and example prompts, then
 121 follow the prompt with the current text sample to estimate the support the sample may have
 122 in data. We clip the log-density-ratio at δ to avoid drawing only training samples.

123 Combining these we get the following maximization objective

$$\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\pi} \left[r_m(x, a) - \beta \log \frac{\pi(a|x)}{\pi_{\text{ref}}(a|x)} - \beta \frac{\pi(a|x)}{\pi_{\text{ref}}(a|x)} + \lambda \log_+ \left(\frac{P(a|x)}{P(\mathcal{D}|x)}, \delta \right) \right] \quad (1)$$

124 where $r_m(x, a) = \frac{1}{E} \sum_{r \in E} r(x, a)$ is the ensemble reward, P is given by GPT probabilities,
 125 $P(\mathcal{D}|x) = P(a^+|x) + P(a^-|x)$, \log_+ is the clipped log function, and β, λ are hyperparameters.

126 4 Experiments

127 **Datasets** We experiment with two public datasets obtained from real-life A/B testing scenarios.

- 128 • **Upworthy:** This dataset records a sample of A/B tests conducted by Upworthy Matias
 129 et al. (2021). The data consists of several versions of headlines created by an editorial
 130 teams for various articles. We only considered text only content and restricted to those
 131 treatments/headlines which were assessed to have statistically different CTRs (at $p=0.10$).
- 132 • **Tweet Popularity:** Tan et al. (2014) studied the effect of wording of a statement on retweeting.
 133 Their tweet popularity dataset is similar to an A/B test with a total of 13k tweet pairs, which
 134 are matched by the topic and the author; where the positive sample is considered as the one
 135 to receive more retweets. We apply a similar pre-processing as Upworthy.

136 **Evaluation** Assessing performance of direct models is straightforward, we simply analyse whether
 137 the model correctly classified test set examples. Evaluating the generative model is more nuanced.
 138 The generative method is assessed based on whether the likelihood of the positive answer is higher
 139 than of the negative answer i.e. whether $\pi(a^+ | x) > \pi(a^- | x)$.

141 Table 1: Accuracy for different approaches on the different A/B
 142 testing datasets. † represents generative models evaluated on
 143 better treatment’s likelihood. * denotes that the prompts triggered
 144 a safety check which were ignored in accuracy calculation

Model	Upworthy	Tweet
GPT-4 PromptOnly	56.6	47.1 *
Claude PromptOnly	58.1	49.6*
Llama-3-8b PromptOnly	55.7	45.3
GPT-4 ICL	64.2	61.3*
Claude ICL	60.1	56.8*
Llama-3-8b ICL	60.7	58.5
OpenAI text-embedding-3-large	82.5	79.9*
Llama-3-8b embedding	74.0	76.5
FT Single	82.8	79.4
FT Ensemble	83.6	80.2
DPO	72.8†	76.1†
Ours	84.5†	81.6†

158 However, the best performance was obtained by fine-tuning these models; Fine-tuned Llama models
 161 outperforms GPT embedding models. These results are consistent between both the Upworthy and

Results Our results are presented in Table 1. From these results we can see that prompting based methods, both direct prompting and few-shot learning, are just a little better than average guessing. Specifically we see no model better than ($\sim 65\%$) accuracy. We also found that giving a few examples for ICL leads to slightly better performance than pure prompt based method ($\sim 60\%$ vs $\sim 55\%$).

Next we also see that methods based on training a model on the data using LLM as representation functions performs much better. With most embedding models we see accuracy of 74% or higher, which is significantly better than prompt based models.

163 Tweet dataset. We also note that the tweet dataset contain samples which triggered safety violations.
 164 With GPT/Claude we excluded these results, and so these numbers are not exactly comparable to
 165 each other.

166 Finally, we tried to rate content by first training an
 167 LLM to produce more engaging content, and using
 168 its likelihood as a measure of rating. The results
 169 indicate that by suitably training the Llama model
 170 to align with the preferences implicitly given by
 171 our dataset, we can match or outperform all the
 172 earlier approaches. This suggests that LLMs can
 173 potentially be used as generators of the treatments
 174 for an A/B test. However since we did not per-
 175 form a human evaluation of its outputs, and more
 176 research is needed in this direction.

177 **Analysis** We delve further into the behaviours
 178 of the different models. As a representative of
 179 the direct evaluation method we chose the GPT
 180 based embedding model, and compared it to the
 181 generative model described earlier. We focus on
 182 the Upworthy dataset here as we have significantly
 183 more number of tests than Tweet. In Figure 3 we
 184 plot the calibration chart i.e. a comparison of
 185 the model accuracy and the predicted probability
 186 of treatment A better than treatment B. For the
 187 generative approach we normalized the probability
 188 of the two considered options instead of using the
 189 output likelihood. The ideal line is of an oracle
 190 calibrated model whose output probabilities will
 191 match its accuracy. From Figure 3 one can see
 192 that the embedding model is overconfident in its
 193 predictions, and modern neural networks are
 194 known to suffer from this (Caruana et al., 2015;
 195 Guo et al., 2017). However, surprisingly the
 196 generative approach seems conservative in its
 197 predictions.

193 We further analyse how the model performance
 194 varies across the difficulty of samples. Difficul-
 195 ty in this context is measured based on the
 196 difference in click-through (CTR) rates. In terms
 197 of downstream impact, having the right decision
 198 when the underlying click rates are different,
 199 is more important than when differences are
 200 lower. In Figure 2, we plot the accuracy of
 201 the GPT3 embedding model against the percentile
 202 of the click rates. We can see that both mod-
 203 els are more accurate as the underlying mean
 204 difference increases. This supports the idea
 205 that the LLM based evaluation can supplement
 206 A/B testing at low risk, as it is less error-
 207 prone when underlying costs of error are higher.

207 5 Conclusion

208 We propose an approach to leverage LLMs for
 209 content experimentation in digital platforms. We
 210 first examined how well LLMs can predict
 211 appealingness of content. First we find that
 212 purely prompt based methods improve over
 213 random chance only by a small factor, suggest-
 214 ing that these methods are not suitable for
 215 predicting engagement. We also try fine-tuning
 216 based approaches to classify content and find
 217 that these are significantly better. Next, we
 218 try to see whether an LLM fine-tuned to pro-
 219 duce engaging content can be used to rate the
 220 treatments. We find that suitably regularized
 221 generative model performs better than the best
 222 fine tuned ensemble models.

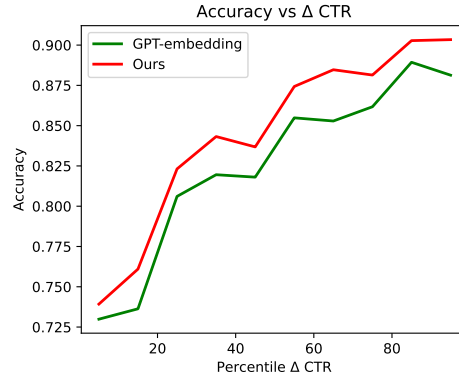


Figure 2: Plot of model accuracy against the absolute value of the mean difference in click rates i.e. $r(a^+, x) - r(a^-, x)$. Both models performs better when the underlying rates are different, but the generative approach outperforms the embedding based model.

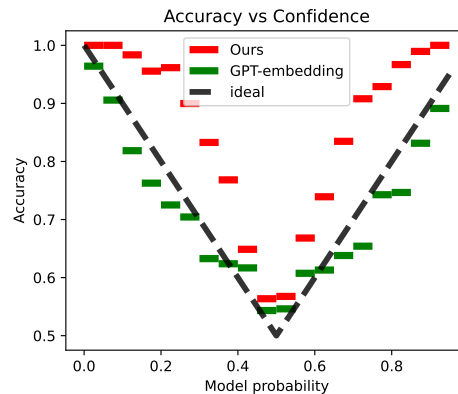


Figure 3: Model accuracy vs model confidence for the generative and embedding approaches. The embedding approach is overconfident while generative model is underconfident in picking the better arm.

217 **References**

- 218 Krisztian Balog and ChengXiang Zhai. 2023. User simulation for evaluating information access
219 systems. In *Proceedings SIGIR*. 302–305.
- 220 Krisztian Balog and ChengXiang Zhai. 2024. User Simulation for Evaluating Information Access
221 Systems on the Web. (2024).
- 222 James Brand, Ayelet Israeli, and Donald Ngwe. 2023. Using gpt for market research. *Available at*
223 *SSRN 4395751* (2023).
- 224 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
225 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models
226 are few-shot learners. *Neurips* (2020).
- 227 Ben Carterette, Evangelos Kanoulas, and Emine Yilmaz. 2011. Simulating simple user behavior for
228 system effectiveness evaluation. In *Proceedings of CIKM*.
- 229 Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015.
230 Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In
231 *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data*
232 *mining*. 1721–1730.
- 233 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024. Self-Play Fine-
234 Tuning Converts Weak Language Models to Strong Language Models. arXiv:2401.01335 [cs.LG]
235 <https://arxiv.org/abs/2401.01335>
- 236 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017.
237 Deep Reinforcement Learning from Human Preferences.
- 238 Anna Coenen. 2019. How The New York Times is Experimenting
239 with Recommendation Algorithms. [https://open.nytimes.com/
240 how-the-new-york-times-is-experimenting-with-recommendation-algorithms-562ff78624d26](https://open.nytimes.com/how-the-new-york-times-is-experimenting-with-recommendation-algorithms-562ff78624d26)
- 241 Thomas Coste, Usman Anwar, Robert Kirk, and David Krueger. 2024. Reward Model Ensembles Help
242 Mitigate Overoptimization. arXiv:2310.02743 [cs.LG] <https://arxiv.org/abs/2310.02743>
- 243 Kawin Ethayarajh. 2019. How Contextual are Contextualized Word Representations? Comparing
244 BERT, ELMo, and GPT-2 Embeddings. In *EMNLP*. 55–65.
- 245 Zhaolin Gao, Jonathan D. Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kianté Brantley,
246 Thorsten Joachims, J. Andrew Bagnell, Jason D. Lee, and Wen Sun. 2024. REBEL: Reinforcement
247 Learning via Regressing Relative Rewards.
- 248 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural
249 networks. In *International conference on machine learning*. PMLR, 1321–1330.
- 250 Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre
251 Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, Johan Ferret, and Mathieu Blondel. 2024. Direct
252 Language Model Alignment from Online AI Feedback. arXiv:2402.04792 [cs.AI] <https://arxiv.org/abs/2402.04792>
253
- 254 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
255 and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International*
256 *Conference on Learning Representations*.
- 257 Junsol Kim and Byungkyu Lee. 2023. AI-Augmented Surveys: Leveraging Large Language Models
258 and Surveys for Opinion Prediction. *arXiv preprint arXiv:2305.09620* (2023).
- 259 Peiyao Li, Noah Castelo, Zsolt Katona, and Miklos Sarvary. 2024. Frontiers: Determining the Validity
260 of Large Language Models for Automated Perceptual Analysis. *Marketing Science* (2024).
- 261 J Nathan Matias, Kevin Munger, Marianne Aubin Le Quere, and Charles Ebersole. 2021. The
262 Upworthy Research Archive, a time series of 32,487 experiments in US media. *Scientific Data* 8, 1
263 (2021), 195.

- 264 Javed Mostafa, Snehasis Mukhopadhyay, and Mathew Palakal. 2003. Simulation studies of differ-
265 ent dimensions of users' interests and their impact on user modeling and information filtering.
266 *Information Retrieval* 6 (2003), 199–223.
- 267 OpenAI. 2024. New embedding models and API updates. [https://openai.com/index/
268 new-embedding-models-and-api-updates/](https://openai.com/index/new-embedding-models-and-api-updates/) Accessed: 2024-05-28.
- 269 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Alex Ray, John Schulman, Jacob Hilton, Fraser
270 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,
271 and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.
272 arXiv:2203.02155
- 273 Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, and Michael S Bernstein.
274 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of ACM USIT*.
- 275 Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun
276 Xu, and Ji-Rong Wen. 2024. Tool Learning with Large Language Models: A Survey.
277 arXiv:2405.17935 [cs.CL] <https://arxiv.org/abs/2405.17935>
- 278 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea
279 Finn. 2024. Direct preference optimization: your language model is secretly a reward model (*NIPS*
280 '23).
- 281 Yuda Song, Gokul Swamy, Aarti Singh, J. Andrew Bagnell, and Wen Sun. 2024. The Importance of
282 Online Data: Understanding Preference Fine-tuning via Coverage. arXiv:2406.01462
- 283 Chenhao Tan, Lillian Lee, and Bo Pang. 2014. The effect of wording on message propagation: Topic-
284 and author-controlled natural experiments on Twitter. arXiv:1405.1438 [https://arxiv.org/
285 abs/1405.1438](https://arxiv.org/abs/1405.1438)
- 286 Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton
287 Murray, and Young Jin Kim. 2024. Contrastive Preference Optimization: Pushing the Boundaries
288 of LLM Performance in Machine Translation. arXiv:2401.08417 [cs.CL] [https://arxiv.org/
289 abs/2401.08417](https://arxiv.org/abs/2401.08417)
- 290 Chenlu Ye, Wei Xiong, Yuheng Zhang, and Tong Zhang. 2024a. Online Iterative RL from Human
291 Feedback with General Preference Model. arXiv:2402.07314
- 292 Zikun Ye, Hema Yoganarasimhan, and Yufeng Zheng. 2024b. LOLA: LLM-Assisted Online Learning
293 Algorithm for Content Experiments. arXiv:2406.02611
- 294 Yangqiaoyu Zhou, Haokun Liu, Tejes Srivastava, Hongyuan Mei, and Chenhao Tan. 2024. Hypothesis
295 Generation with Large Language Models. arXiv:2404.04326 [cs.AI] [https://arxiv.org/abs/
296 2404.04326](https://arxiv.org/abs/2404.04326)
- 297 Lixi Zhu, Xiaowen Huang, and Jitao Sang. 2024a. How Reliable is Your Simulator? Analysis
298 on the Limitations of Current LLM-based User Simulators for Conversational Recommendation.
299 arXiv:2403.16416
- 300 Yutao Zhu, Huaying Yuan, Shuting Wang, Haonan Chen, Zhicheng Dou, and Ji-Rong Wen. 2024b.
301 Large Language Models for Information Retrieval. arXiv:2308.07107 [cs.CL] [https://arxiv.
302 org/abs/2308.07107](https://arxiv.org/abs/2308.07107)
- 303 Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. [n. d.].
304 Can large language models transform computational social science? *Computational Linguistics*
305 ([n. d.]).

306 **Appendix**

307 The prompt used for in-context learning is given in Example 1. Instructions from this are also used in
308 the context when tuning the models.

```
309 

---


310 Instruction Prompt
311 You are an expert marketing writer for a newspaper company. You you are excellent at choosing which
312 headlines are likely to get more clicks for an article.
313 You will be given an article context and two headlines, from which you determine which headline was
314 clicked more often.
315 You are given the headlines as "Headline _" where _ is either 1 or 2. Give your final answer in the
316 following format:
317 "Answer: Headline _"
318
319 User Prompt
320 Here are some previous examples to help you:
321 ... more examples here ...
322 Which of the following headlines has more clicks:
323 Article: <context>
324 Headline 1: <headline_1>
325 Headline 2: <headline_2>
326
327 Think step by step, and explain your reasons
328 Step 1: Look at the new pair of headlines and compare them with the examples associated with each
329 pattern.Step 2: Find the set of examples that is closest to the given pair of headlines, and pick the
330 pattern associated with that set of examples.
331 Step 13: Think about which one out of the pair of headlines will get more clicks.
332 Step 14: Give your final answer.
333 

---


```

Example 1: Zero/Few-shot Inference.

334 We further analyse how the model performance varies across the the statistical significance of
335 the difference in the CTR rates ¹. The significance score is done by a Welch-t test (two-sample
336 uncommon variance t-test).

337 This is different from Figure 2 as difficulty in
338 this context is measured based on the statisti-
339 cal significance of the difference in the CTR
340 rates. A high difference in click through rates
341 need not mean high significance, as the dif-
342 ference is adjusted for the variance and/or the
343 number of impressions for computing the sig-
344 nificance. Note that since we already filtered
345 out non-conclusive tests, we are considering
346 only low p-value samples. The result of accu-
347 racy on the test-set for is presented in Figure
348 4. We can see that both models in general are
349 more accurate as the significance increases
350 (p-value decreases).

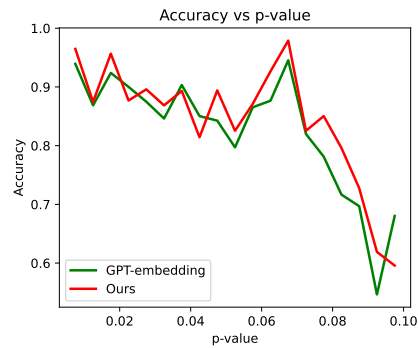


Figure 4: Plot of model accuracy against the per- centiles of the click-rates difference of different arms.

¹Since the significance score is also dependent on number of impressions which get influenced by the experimenter decisions, the difference in rates is not an ideal measure of difficulty.