

Improving Generative Ad Text on Facebook using Reinforcement Learning

Daniel R. Jiang^{†,*}, Alex Nikulkov[†], Yu-Chia Chen, Yang Bai, Zheqing Zhu

Meta Platforms, Menlo Park, California, USA.

[†]These authors contributed equally to this work.

*Corresponding author. Email: drjiang@meta.com

Abstract

Generative artificial intelligence (AI), in particular large language models (LLMs), is poised to drive transformative economic change. LLMs are *pre-trained* on vast text data to learn general language patterns, but a subsequent *post-training* phase is critical to align them for specific real-world tasks. Reinforcement learning (RL) is the leading post-training technique, yet its economic impact remains largely underexplored and unquantified. We examine this question through the lens of the first deployment of an RL-trained LLM for generative advertising on Facebook. Integrated into Meta’s *Text Generation* feature, our model, “AdLlama,” powers an AI tool that helps advertisers create new variations of human-written ad text. To train this model, we introduce *reinforcement learning with performance feedback* (RLPF), a post-training method that uses historical ad performance data as a reward signal. In a large-scale 10-week A/B test on Facebook spanning nearly 35,000 advertisers and 640,000 ad variations, we find that AdLlama improves click-through rates by 6.7% ($p = 0.0296$) compared to a supervised imitation model trained on curated ads. This represents a substantial improvement in advertiser return on investment on Facebook. We also find that advertisers who used AdLlama generated more ad variations, indicating higher satisfaction with the model’s outputs. To our

knowledge, this is the largest study to date on the use of generative AI in an ecologically valid setting, offering an important data point quantifying the tangible impact of RL post-training. Furthermore, the results show that RLPF is a promising and generalizable approach for *metric-driven post-training* that bridges the gap between highly capable language models and tangible outcomes.

1 Introduction

Generative artificial intelligence (AI) is increasingly being recognized for its transformative potential across industries, driving innovations in content creation [1, 2, 3, 4, 5], education [6, 7, 8], medicine [9, 10, 11, 12], and complex decision-making [13, 14, 15]. It is also widely believed to have the potential for significant economic impact [16, 17, 18, 19, 20], with pioneering work showing quantifiable benefits in the areas of customer support [21, 22], enterprise information tasks [23], legal tasks [24, 25], consulting work [26], software development [27, 28, 18], email marketing [29], online advertising [30], and a variety of professional writing tasks [31]. However, turning this potential into consistent real-world impact depends critically on how models are fine-tuned and aligned during a phase of model training called *post-training*.

The initial phase for large language models (LLM) training is called *pre-training*, where the model learns general language patterns and world knowledge from large-scale unlabeled text data. However, to perform effectively in real-world applications, generative AI systems require some form of *post-training*—a suite of methods designed to adapt pre-trained foundation models to specific downstream tasks [32]. A common approach involves supervised fine-tuning (SFT), in which models are trained to follow human instructions by mimicking curated target responses [33, 34]. Another example is the use of *reinforcement learning from human feedback* (RLHF) to train the model on human preference data [35], which has been shown to be highly effective in domains like summarization [36], chatbot assistants [37], and instruction-following [38]. More recently, researchers have found that reinforcement learning (RL) using *verifiable rewards* (i.e., objective signals from domains like math or coding that can be automatically validated) can lead to marked

improvements in LLM reasoning and problem-solving capabilities [32, 39].

As noted above, there is a broadening body of research focused on quantifying the impact of LLMs across various sectors. However, the specific impact of *the post-training phase* has been relatively underexplored, despite its critical role in the development of virtually all prominent LLM models. In this paper, we present new findings on the concrete impact of RL for post-training. Our analysis is conducted through the perspective of the online advertising industry, which is a major driver of the global economy: worldwide expenditure in online ads is projected to reach \$513 billion USD in 2025 [40] and is expected to represent 63% of the total global advertising revenue across all mediums [41].

Specifically, we focus on Meta’s Text Generation product, which uses an LLM to generate variations of an advertiser’s human-written ad text. This allows each advertiser to select multiple versions of the ad to show to potential customers, taking advantage of Meta’s ad delivery systems to show the most performant variants. The initial version of the Text Generation product used an LLM that was trained in a naive manner, by fine-tuning it to imitate the style of a set of curated ads using supervised learning (SFT).

In this paper, we take the next step and describe our efforts in using reinforcement learning to improve the Text Generation LLM, with the measurable goal of writing *more engaging* ad text. We do this by directly using performance (i.e., click-through rates) as a reward signal. We call the approach *reinforcement learning with performance feedback* (RLPF). By viewing each user’s behavior (i.e., click or no click) on each ad impression as a small piece of human feedback, RLPF can be thought of as an extension of the well-known RLHF technique [35, 33], except that each piece of ad text is being “rated” by *thousands of humans* (i.e., Facebook users who see an ad impression) via click or no click signals. In that sense, one can view RLPF as falling somewhere in between traditional RLHF and the recently explored direction of RL with verifiable rewards [32].

We report on the results of a large-scale online experiment (i.e., A/B test) conducted on Facebook that compared the effect of providing advertisers with the RLPF-trained model versus the naive SFT-based imitation model. Our experiment and analysis encompass approximately 35,000 advertisers

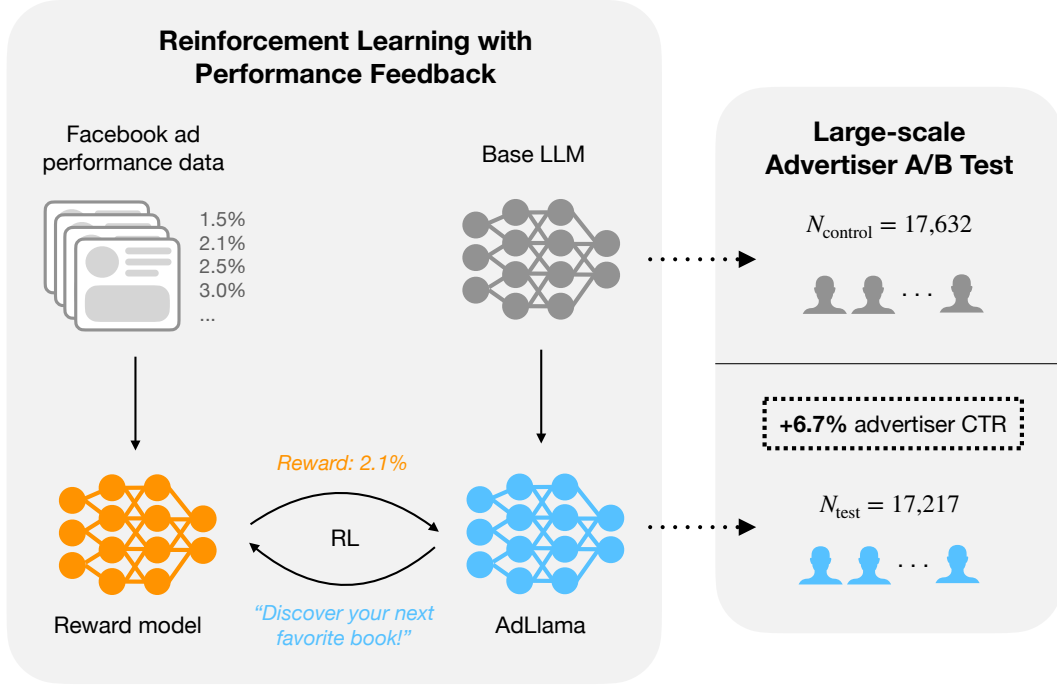


Figure 1: **Overview of our contributions.** Our first contribution, illustrated in the left panel, is RLPF, a reinforcement learning (RL) approach to post-training an LLM based on an aggregate performance metric. We apply RLPF to a generative AI feature in Meta’s Ads Manager that helps advertisers generate new ad text variations. To do this, we use Facebook ad performance data (i.e., click-through rates) to train a reward model, which can score the effectiveness of a piece of ad text. Subsequently, we align the LLM toward this reward model using RL, which involves iteratively generating and scoring ad text, illustrated in lower part of the figure. The resulting LLM is called “AdLlama.” Our second contribution, illustrated in the right panel, are the results of a large-scale advertiser A/B test, encompassing approximately 35,000 advertisers and 640,000 ad variations, which showed that advertisers who received AdLlama achieved 6.7% higher advertiser-level click-through rates (CTR). To our knowledge, this study is the *largest reported so far* that investigates the use of generative AI in an ecologically valid setting.

and 640,000 ad variations, observed over 10 weeks. The results indicate that when equipped with the RLPF model, advertisers achieved a click-through rate (CTR) increase of 6.7% compared to the naive imitation model. In addition, the number of ad variations created by the advertiser increased by 18.5%, suggesting that advertisers were more willing to adopt AI suggestions when they came from AdLlama.

The implications of this result are significant. First, this improvement highlights the effectiveness of reinforcement learning as a methodology for metric-driven post-training of LLMs,

especially in business use cases. Second, it offers a valuable data point by quantifying the benefits of RL-based post-training of LLMs in the domain of online advertising, contributing an important insight in the broader ongoing effort to understand the implications of generative AI.

To our knowledge, is the *largest study to date* examining the use of generative AI in a real-world, ecologically valid context. Previous analyses often relied on smaller-scale experiments involving a few thousand participants [26, 27, 22, 28, 29] or controlled laboratory settings using crowd workers [31, 24, 25, 30, 20]. These limitations largely reflect the high computational costs of training and deploying such models, as well as the scarcity of opportunities for large-scale field deployment.

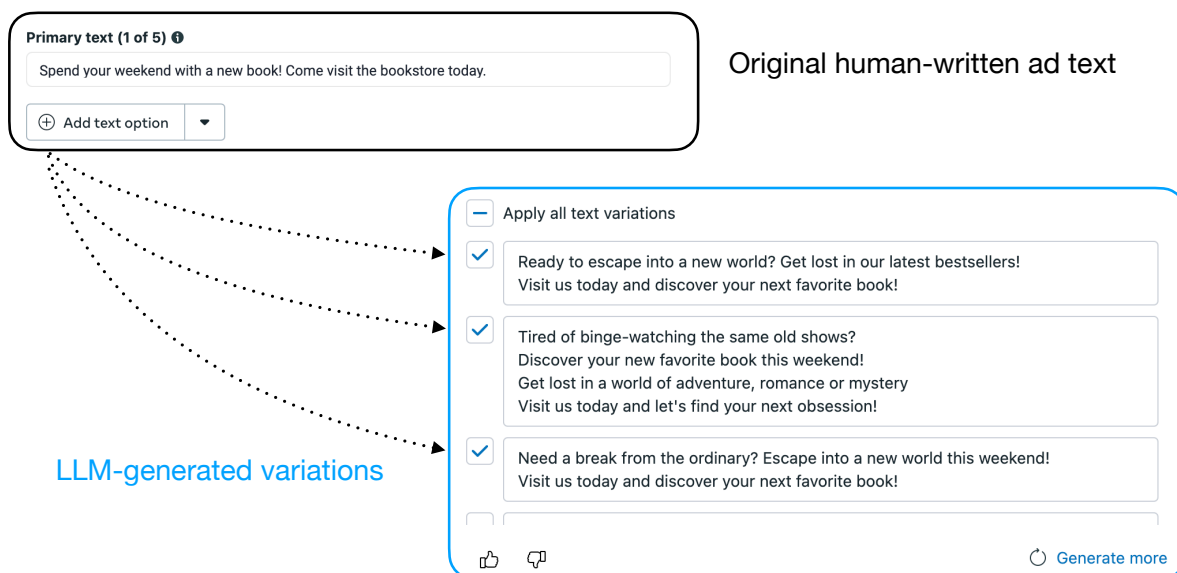


Figure 2: **The user interface of the Meta Text Generation product.** Here we show an example of how the Text Generation product functions during the ad creation process. The advertiser provides an original version of the ad text (“*Spend your weekend with a new book! Come visit the bookstore today.*”), which serves as input to Text Generation. The underlying LLM then produces a set of ad text variations, allowing the advertiser to pick and choose the ones they prefer to use. The advertiser can also opt to generate additional variants via the “Generate more” button.

2 Meta’s Text Generation Product

The Text Generation product [42] is a generative AI feature within Meta Ads Manager that allows advertisers to experiment with multiple versions of their ad text. The feature works by taking an

advertiser’s original ad text as input and then using an underlying LLM to suggest new variations, which may, for example, emphasize key selling points or add additional creative messaging.

2.1 User Interface

The user interface of the Text Generation product is illustrated in Figure 2. The flow is as follows. Advertisers first input an original ad text. The underlying LLM then generates and displays multiple text variations. Advertisers can then select the variations they want to use, edit them directly in the text box, or even add their own custom variations via the “Add text option” button [43]. In Figure 2, we show an example where the advertiser has selected a total of *four* text variations, the original ad text along with three AI-generated variations. The option to continue generating additional text variations for consideration is available via the “Generate more” button, but the advertiser is limited to selecting a maximum of five variations for delivery to users.

We note that although the Text Generation interface appears during the ad creation process, advertisers are *never required* to select AI-written ads. Some other possible actions that the advertiser may opt to take are: (1) ignore the AI’s suggestions completely and continue with the original text, (2) ignore the suggestions and instead supply multiple human-written text variations, (3) edit the AI-written ads before selecting them, (4) use the AI-written suggestions as inspiration for additional human-written variations, or (5) make AI-inspired edits to the original human-written variation after seeing the AI-generated variants. Therefore, the LLM can play an nuanced role in shaping ad text, regardless of whether the advertiser ultimately selects the precise wording that was generated.

2.2 Naive Imitation Models

The original (naive) version of Text Generation LLM was released to advertisers in November 2023. This LLM is based on Meta’s open-source foundation language model, the 7 billion (7B) parameter version of Llama 2 Chat [44]. As mentioned previously, the Llama model was post-trained to imitate a set of curated ads using SFT. We refer to this model as “Imitation LLM v1.” We refer to an

improved version of the imitation model that uses higher quality data as “Imitation LLM v2.” The main difference between Imitation LLM v1 and v2 is that while the v1 training dataset is fully based on synthetic generations from a larger LLM, the v2 data additionally included human-written (i.e., contractor-written) examples. These training examples, whether synthetic or human-written, were curated by asking either the LLM or human to rewrite existing ads using specific instructions, such as “paraphrase and shorten,” “make clear,” “make actionable,” “empathize,” “pose as a question,” or “focus selling point.” We detail the training process of both Imitation LLM models in Section A.4 of the supplementary materials.

The work presented in this paper is subsequent to the initial release of the Text Generation product. Our goal in this paper is to improve the initial imitation-based Text Generation LLM, focusing on *quantifiably improving advertiser performance* relative to the original model, as measured in terms of click-through rate (CTR). We tackle this problem using a new idea, reinforcement learning on aggregate performance feedback signals.

3 Methods

In this section, we first describe the methodology (including data preparation, reward model design, and reinforcement learning) for training the new version of the Text Generation LLM, which we call “AdLlama.” We then discuss the design of the experiment (i.e., A/B test) used in quantifying the performance improvement of the new model compared to the naive imitation model.

3.1 Reinforcement Learning with Performance Feedback (RLPF)

Although pre-trained LLMs have absorbed vast amounts of knowledge, they are typically not ready for widespread use. A crucial step before deployment to users is to *align* the model for use on downstream tasks, which may include summarization [36], answering questions [45], engaging in dialogue [46], or following a wide range of instructions from users [33]. The primary approach for LLM alignment involves collecting *preference* data from human labelers, who compare two

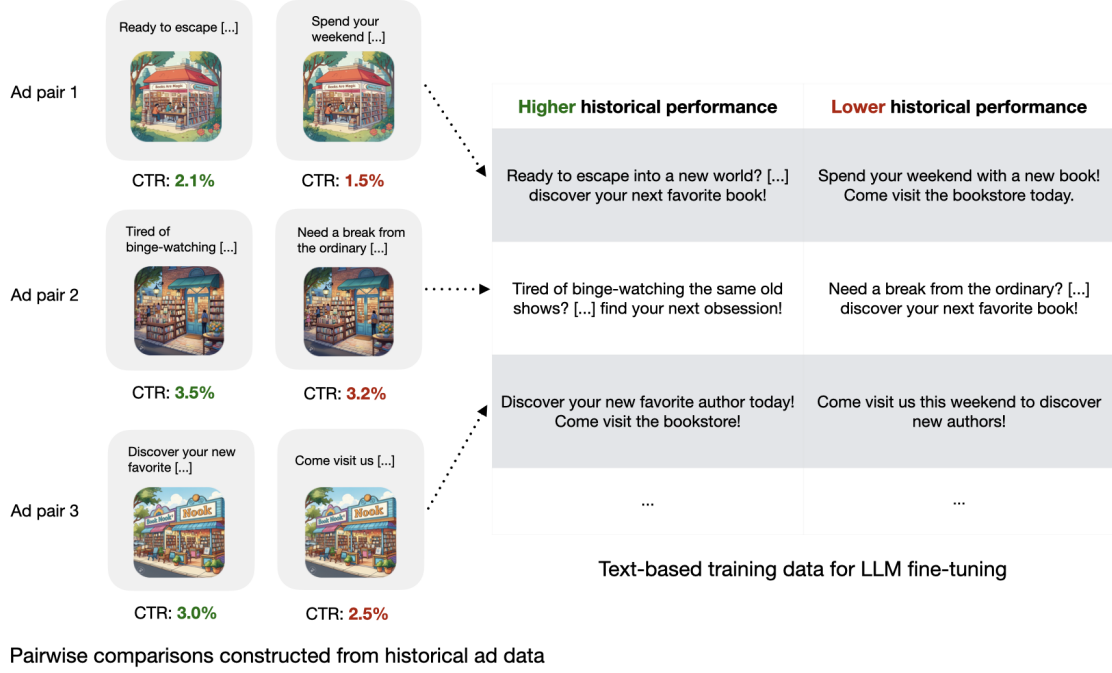


Figure 3: **Pairwise training data from historical multitext performance outcomes.** Based on *multitext data*, data where advertisers change only the ad text while keeping all other ad components fixed, we are able to construct a pairwise training data point that distinguishes between text variations that are preferred and not preferred, as measured by CTR. On the left, we show historical ad pairs that have different text, but the same image. On the right, we show that each ad pair contributes to a single preference row for our training data.

responses and indicate which is better. The model is then fine-tuned on the preference data, encouraging it to generate responses more similar to those preferred by humans [35, 33]. This process is known as *reinforcement learning with human feedback* (RLHF). Because the “quality” of LLM responses for many standard LLM tasks (e.g., open-ended dialogue, creative writing) is subjective and hard to quantify in nature, training on human preferences is often the closest that one can come to a well-defined optimization objective.

Our key insight is that, unlike the above LLM use-cases, the task of *crafting effective ad text* can be clearly associated with a *quantifiable and measurable* objective: the CTR¹ of the ad. Note that such a setup holds anytime there is a concrete performance metric and exists beyond online

¹The CTR is defined as the ratio of number of ad clicks to the total number of ad impressions (i.e., views). We note that CTR is typically considered a proxy metric for true goal of *conversion rate* (the number of conversions, or purchases, per impression). However, because conversion rate is a far noisier metric to work with due to the sparsity of conversions, we resort to CTR.

advertising (e.g., e-commerce, AI customer support agents, ed tech). We propose the following *general* approach, which can be thought of as a metric-driven extension of RLHF:

1. First, using aggregate performance metrics, we train a performance *reward model*, a model that can assign a reward (i.e., a score) to a piece of text, where more performant text are given higher scores.
2. Subsequently, we use the reward model as an interactive environment to perform RL fine-tuning, where the goal is to fine-tune the LLM to be more likely to generate high-reward text.

We now describe how we train a CTR-based performance reward model. Even prior to the launch of the Text Generation product, there existed a common practice among advertisers of testing multiple (human-written) text variants for a single ad using one of Meta’s advertiser tools called “Multiple Text Optimization” [47]. This practice enables us to observe historical ad data where, except for the text, all other components of the ad—such as the image, title, and targeting criteria—remain constant. We refer to this as *multitext* data.

From the multitext data, we are able to construct *preference pairs*, where the higher CTR text is marked as “more preferred” and the lower CTR is marked as “less preferred.” See Figure 3 for an illustration of this process. We call this the *pairwise* dataset, which supports the standard Bradley-Terry preference-based approach to reward model training [48, 35]. We also considered a more naive reward modeling approach via a *pointwise* dataset, where each row is simply the ad text and its resulting CTR. We can then use a standard supervised learning to train a reward model directly using CTR as labels. However, we found that the pointwise reward model was less capable of discerning ordering (or ranks) between similar pieces of ad text, which is ultimately more important than purely predicting CTR [49, 50]. Further details on the data curation and reward model training are in Sections A.1 and A.2 of the supplementary materials.

Given a trained reward model r_θ , we use the proximal policy optimization (PPO) algorithm [33, 51] to align the LLM with high-performance ad text. We added a length penalty to counteract

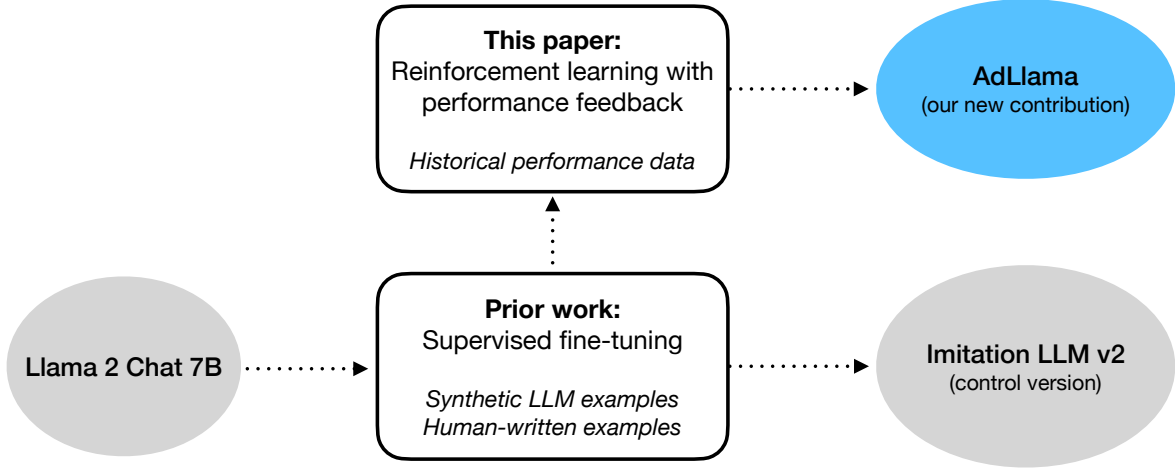


Figure 4: **AdLlama versus Imitation LLM v2.** Both AdLlama and Imitation LLM v2 originate from a base 7B Llama 2 Chat model. The difference is that AdLlama is further trained via RLPF and historical ad performance data, while Imitation LLM v2 is only trained using SFT to imitate a set of curated examples (includes both LLM-generated synthetic examples and human-written examples).

the tendency for the model to generate ad text that is too long, leading to the following PPO optimization formulation:

$$\max_{\pi_{\phi}} \mathbb{E}_{x \sim \mathcal{D}_{\text{LLM}}, y \sim \pi_{\phi}(y|x)} [r_{\theta}(x, y) - \beta \text{KL}[\pi_{\phi}(\cdot | x), \pi_{\text{ref}}(\cdot | x)] - \alpha \text{length}(y)],$$

where π_{ϕ} is the LLM, \mathcal{D}_{LLM} is the pointwise dataset used for LLM post-training, $\text{KL}(p, q)$ is the Kullback-Leibler divergence between two probability distributions p and q , $\text{length}(y)$ is the number of tokens in y , and β and α are weight parameters.

Full details of the approach are given in Section A.2 of the supplementary materials. We used the RLPF technique to improve Imitation LLM v2, which, like Imitation LLM v1, is based on the 7B version of Llama 2 Chat [46]. We refer to our RLPF-based ad text generation model as “AdLlama.” An illustrative comparison of the models is given in Figure 4, where we emphasize both the difference in training method (RLPF versus SFT) and the training data (historical ad performance versus curated examples).

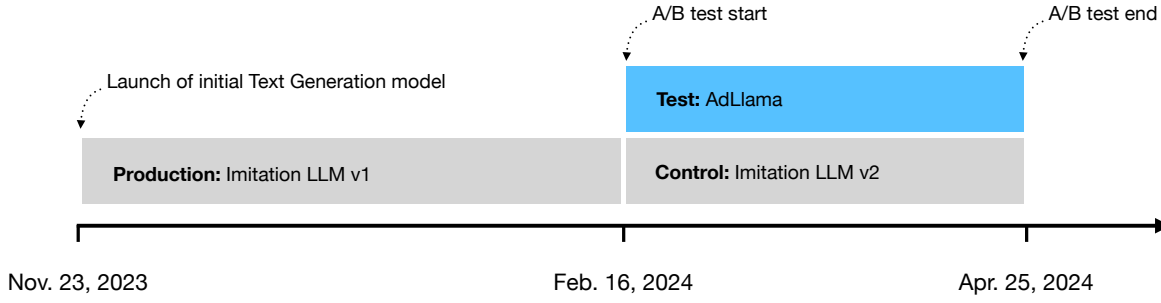


Figure 5: **A/B test timeline.** Our A/B test ran from February 16, 2024 until April 25, 2024. Prior to that, Imitation LLM v1 was launched on November 23, 2023 as the initial version of the Text Generation LLM.

3.2 Experiment Design

We conducted a large-scale A/B test (i.e., randomized control trial) to evaluate the impact of RLPF training on advertiser performance by comparing the AdLlama model against Imitation LLM v2. The A/B test ran for 10 weeks, from February 16, 2024 until April 25, 2024, on $N = 34,849$ advertisers based in the United States. We randomize at the advertiser level: each advertiser is randomly assigned either (1) the Imitation LLM v2 (“control”) or (2) the AdLlama LLM trained with RLPF (“test”). Figure 5 shows the timeline of the A/B test and how it relates to the launch of the initial Imitation LLM v1 model.

Our primary focus is on *advertiser-level performance*, defined as performance aggregated over all *direct-response*² ads created by the advertiser³ over the 10-week experimental period. We chose to examine advertiser-level performance because our goal is to understand how generative AI can improve advertiser return on investment. More specifically, our analysis centers on the following metrics: total engagement (clicks), total impressions (views), total number of ads created, and total

²There are two broad categories of ads, *direct-response* and *brand/awareness*. *Direct-response* ads are those that have a goal of optimizing for some form of engagement or immediate action from the user. Examples of direct-response advertising goals include clicks, purchases, app installs, sign ups, donations, or subscriptions. On the other hand, *brand* ads aim to drive awareness simply through views of the ad (or video), and the goal is not to drive engagement from the user. Since the AdLlama model is designed to improve engagement, the relevant metric for evaluating AdLlama is CTR of direct-response ads.

³A frequent practice of advertisers is to duplicate existing ads, some of which may have been originally created prior to the experiment start date. Such ads are excluded from our results, but duplicates of ads created *during* the experimental period are included.

number of ad variations created,⁴ all of which are defined at the advertiser level across the 10-week experimental period on Facebook mobile feed.

We also log a number of advertiser covariates, which are listed in full in Section A.3.1 of the supplementary materials, along with detailed descriptions. Two notable covariates are: the advertiser’s pre-experiment lifetime CTR across all Meta apps (“pre_exp_ctr”) and whether or not the advertiser is “new,” defined as whether the advertiser’s pre-experiment ads have accumulated fewer than 1,000 impressions (“is_new_advertiser”). We also include the advertiser’s ad creation behavior during the period after Imitation LLM v1 was released, but before the experiment started (“nov_feb_ad_cnt”, “nov_feb_variant_cnt”, “has_created_llm_ad”). This period of time is highlighted in blue in Figure 5 and represents the initial few months of the Text Generation feature being available to advertisers.

We also include other advertiser characteristics: vertical, expertise level, budget level, business account status, and account age. In addition to covariate details, we also present descriptive statistics and the balance of these covariates across conditions in Section A.3.1 of the supplementary materials (see Table A.3.1). We do not observe statistically significant imbalances between the two groups.

4 Main Results

4.1 Advertiser Performance

We aim to evaluate the effect of AdLlama and Imitation LLM v2 on advertiser-level CTRs. Our main regression specification is a log-binomial⁵ regression model, with the left-hand side representing

⁴Recall from Figure 2 that in the Text Generation product, a single “ad” is associated with one or more “variations” where the body text is altered.

⁵Binomial regression is a natural candidate for modeling a CTR, since we are interested in clicks (successes) relative to impressions (total trials). Our use of a log link function leads to a relative risk rather than an odds ratio [52]. The relative risk is more desirable in our setting since it can be directly interpretable as a ratio of CTRs.

Table 1: Log-binomial regression results on advertiser-level (engagement, impressions), with the “treatment” variable being the indicator for using AdLlama. The variable “log_pre_exp_ctr_existing” refers to the product $\mathbf{1}_{\text{existing}(i)} \cdot \log(\text{CTR}_i^{\text{pre}})$. Coefficients are reported on the link function scale (log scale); the 6.7% improvement quoted in the main text computed by $6.7\% \approx \exp(0.0651) - 1$. Fixed effect terms are omitted for brevity. All reported p-values are two-sided.

	<i>Dependent variable:</i>
	engagement/impressions
treatment	0.0651** (0.0299)
log_pre_exp_ctr_existing	0.5931*** (0.0264)
is_new_advertiser	−2.3469*** (0.1569)
pre_exp_ad_cnt	0.0024 (0.0028)
pre_exp_impressions	−0.0001 (0.0001)
pre_exp_engagement	0.0036 (0.0049)
account_age_yr	0.0008 (0.0038)
nov_feb_ad_cnt	0.0013 (0.0019)
nov_feb_variant_cnt	−0.0003 (0.0006)
is_business_account	−0.0286 (0.0424)
has_created_llm_ad	−0.0109 (0.0345)
Constant	−1.1233*** (0.1506)
Observations	34,849

Note: HC1 robust standard errors in parentheses. *p<0.1; **p<0.05; ***p<0.01

the logarithm of the CTR of advertiser i :

$$\begin{aligned} \log \mathbb{E}(Y_i/n_i | X_i) = & \beta_0 + \beta_1 \cdot \text{AdLlama}_i + \beta_2 \cdot \mathbf{1}_{\text{existing}(i)} \cdot \log(\text{CTR}_i^{\text{pre}}) + \beta_3 \cdot \mathbf{1}_{\text{new}(i)} \\ & + \beta_{4:10}^T Z_i + \theta_{\text{vert}(i)} + \alpha_{\text{budget}(i)} + \lambda_{\text{expertise}(i)}, \end{aligned} \quad (1)$$

where Y_i is the number of clicks (engagement) for the i -th advertiser, n_i is the impression count, X_i is the set of all covariates, $\text{CTR}_i^{\text{pre}}$ is the pre-experiment CTR (“pre_exp_ctr”), AdLlama_i is a binary indicator for the AdLlama LLM treatment (which is randomly assigned), $\mathbf{1}_{\text{existing}_i}$ is an indicator for “is_new_advertiser equal to 0,” $\mathbf{1}_{\text{new}_i}$ is an indicator for “is_new_advertiser equal to 1,” $\theta_{\text{vert}(i)}$ are vertical fixed effects, $\alpha_{\text{budget}(i)}$ are budget category fixed effects, $\lambda_{\text{expertise}(i)}$ are expertise category fixed effects, and Z_i is the vector of remaining numerical covariates. Note that $\mathbb{E}(Y_i/n_i | X_i)$ is the CTR of advertiser i .

Because some advertisers in the experiment had no pre-experiment ad impressions (or an insufficient number of impressions for a reliable CTR calculation), it is not possible to naively

include the pre-experiment CTR as a covariate for all advertisers. Instead, we devise the following strategy to properly account for pre-experiment CTR through the β_2 and β_3 terms of Equation 1. If the advertiser is an “existing” advertiser (i.e., has sufficient ad impressions), then we directly incorporate its pre-experiment CTR via the term $\beta_2 \cdot \mathbf{1}_{\text{existing}(i)} \cdot \log \text{pre_exp_ctr}_i$. This can be interpreted as a baseline log CTR for that advertiser, which is then “adjusted” by the regression specification based on other covariates. New advertisers do not have a reliable pre-experiment CTR, so we estimate a baseline log CTR using the term $\beta_3 \cdot \mathbf{1}_{\text{new}(i)}$. In Section A.3.2 of the supplementary materials, we show how our regression specification above leads to an intuitive interpretation.

Table 1 presents the findings from the log-binomial regression analysis. The results indicate that AdLlama provides a statistically significant **6.7% increase** in advertiser-level CTR ($p = 0.0296$ and a standard error of 0.0299) when compared to the naive Imitation LLM v2. This roughly corresponds to an absolute increase in advertiser-level CTR from 3.1% to 3.3%. Although the absolute increase appears modest at first glance, a 6.7% relative increase in CTR represents a substantial improvement in an advertiser’s return on investment for Facebook ads. Furthermore, on mature, highly-optimized ad platforms like Facebook, even small increases in CTR are typically difficult to achieve.

4.1.1 Robustness Checks

We provide several robustness checks to validate the findings above. Since we did not observe significant imbalances between the control and test groups, we used model-free results as additional supporting evidence (see Section A.3.3 of the supplementary materials). Further, in Section A.3.4 of the supplementary materials, we tested various alternative CTR regression specifications, including quasi-binomial, logistic, Poisson, and quasi-Poisson regressions. These analyses yielded qualitatively similar effects. Finally, in Section A.3.5 of the supplementary materials, we conducted separate linear regressions on clicks and impressions, providing statistical evidence that AdLlama increases the total number of clicks per advertiser, while it did not affect the total number of impressions delivered for the advertiser. This is further supporting evidence that CTR is increased

under AdLlama.

4.2 Impact on Ad Variations Created

We are also interested in how AdLlama and Imitation LLM v2 affect advertisers’ usage of the Text Generation product. We consider two outcomes, (1) the number of *ad variations* created during the experimental period and (2) the number of ads created during the experimental period. Recall that a single “ad” is associated with one or more “variations” (see Figure 2).

We use a linear regression with the same covariates as we did in the log-binomial regression of Equation 1, except without the log-transformation⁶ of the pre-experimental CTR:

$$\begin{aligned} \text{Outcome}_i = & \beta_0 + \beta_1 \cdot \text{AdLlama}_i + \beta_2 \cdot \mathbf{1}_{\text{existing}(i)} \cdot \text{CTR}_i^{\text{pre}} + \beta_3 \cdot \mathbf{1}_{\text{new}(i)} \\ & + \boldsymbol{\beta}_{4:10}^T \mathbf{Z}_i + \theta_{\text{vert}(i)} + \alpha_{\text{budget}(i)} + \lambda_{\text{expertise}(i)} + \epsilon_i. \end{aligned} \quad (2)$$

Here, Outcome_i refers to either advertiser i ’s variation count or ad count and ϵ_i is the error term.

Table 2 reports the results of these regressions. We observe strong evidence that usage of the AdLlama LLM *increases the number of ad variations created by the advertiser*, while the total number of ads created remained statistically the same. Specifically, we see that the number of ad variations increased by 3.1 ($p < 0.01$) from roughly 16.8 variations (Imitation LLM v2) to 19.9 (AdLlama). This is an 18.5% increase when using the AdLlama LLM. This suggests that for each ad, advertisers were more willing to use the Text Generation product’s suggestions when they came from AdLlama compared to Imitation LLM v2.

5 Discussion

Our work shows that RL with performance feedback can be used to train LLMs to generate ad text that resonates with advertisers and drives measurable engagement from users on Facebook.

⁶We remove the log-transformation because there is no longer a log link function; see Section A.3.2 of the supplementary materials for further discussion.

Table 2: Linear regression results on advertisers’ ad creation behavior: “variant_cnt” is the number of ad variations created, “ad_cnt” is the number of ads created, and “pre_exp_ctr_existing” refers to the term $\mathbf{1}_{\text{existing}(i)} \cdot \text{CTR}_i^{\text{pre}}$. All other notation remains consistent with Table 1. We report HC1 robust standard errors. Fixed effect terms are omitted for brevity. All reported p-values are two-sided.

	<i>Dependent variable:</i>	
	variant_cnt	ad_cnt
treatment	3.113*** (0.528)	0.040 (0.189)
pre_exp_ctr_existing	−5.915 (11.057)	−6.568* (3.549)
pre_exp_ad_cnt	3.803*** (0.689)	1.564*** (0.280)
pre_exp_impressions	0.009 (0.011)	0.002 (0.004)
pre_exp_engagement	−0.333 (0.254)	−0.123 (0.098)
account_age_yr	−0.397*** (0.083)	−0.197*** (0.029)
nov_feb_ad_cnt	0.750* (0.441)	0.615*** (0.170)
nov_feb_variant_cnt	0.096 (0.137)	−0.055 (0.050)
is_new_advertiser	2.656*** (0.776)	0.991*** (0.268)
is_business_account	3.405*** (0.420)	1.416*** (0.148)
has_created_llm_ad	2.078** (1.033)	0.280 (0.399)
Constant	7.797*** (2.035)	3.831*** (0.685)
Observations	34,849	34,849
R ²	0.091	0.116
Adjusted R ²	0.091	0.115

Note: HC1 robust standard errors in parentheses.

*p<0.1; **p<0.05; ***p<0.01

Specifically, our large-scale A/B test on Meta’s Text Generation product shows that the RLPF-based model significantly increases advertiser-level CTRs, along with the number of ad text variations that advertisers were willing to employ. These results support the concept of anchoring the fine-tuning process in real-world, aggregate performance metrics, rather than relying solely on human raters’ preference feedback or rule-based rewards.

5.1 Limitations

There are several limitations of our work that we now discuss. Our model was trained using offline historical performance data. Therefore, this is equivalent to a single round of *offline* RL, where there is no real-time interaction with the environment. To further refine our model, we could incorporate the performance outcomes of LLM-generated ads in an iterative process. This approach would align more closely with *online* RL, where the model continuously interacts (by taking actions) with

the environment and adapts based on real-time feedback in the form of rewards and transitions [53]. Such a system would be more capable of adapting to new trends and perhaps even discovering new ones through *exploration*, the process of experimenting with new actions (i.e., trying out new ad text formats unseen in the historical data).

Our current model primarily focuses on ad performance, but other factors are also important to consider. As an example, there may be a trade-off between generating ads that perform well and those that exhibit high creativity. Additionally, the model’s ability to adhere to specific advertiser instructions, such as maintaining a particular tone, is another important consideration. Addressing these aspects would require a multi-objective optimization approach to balance various objectives effectively. Finally, our model currently does not take into account the human component of the Text Generation product: before an ad text variation can be delivered to users, the advertiser must explicitly select that variation for delivery. An alternative way to train the future iterations of the RLPF reward model is to weigh the CTR by the likelihood that the text is selected by advertisers.

Beyond individual ad performance, platform-level factors, such as the diversity of the ad inventory, are also important for a positive user experience. Future work should explore strategies that simultaneously consider these other factors, while also optimizing for performance.

5.2 Broader Implications

Our findings contribute to the growing body of literature on understanding the impact of LLMs. By quantifying the benefits of RL-based post-training in online advertising, we provide a concrete data point that highlights the potential for these models to ingest relevant performance metrics and subsequently create real business impact. The ability to generate more engaging ad content not only improves existing advertisers’ return on investment, but could also lower the barrier to entry for new and inexperienced advertisers (e.g., small businesses) by reducing the need for extensive marketing expertise and resources.

Our methodology is not limited to online advertising: the principles of RLPF can be adapted to other domains where aggregate performance metrics are available. By using performance data as a

feedback mechanism, organizations can fine-tune LLMs to optimize for their desired outcomes. For example, the core methodology can easily be extended to closely related settings like personalized email campaigns or e-commerce product descriptions. RLPF can also be extended to settings with multiple rounds of interactive feedback, such as AI customer support agents, using metrics like resolution rates, satisfaction scores, or user response times.

There are also less obvious settings where RLPF could be applied. For example, in online learning platforms, student performance data (test scores and engagement metrics) could guide the generation of adaptive learning content, while for certain public awareness campaigns (e.g., vaccination, energy consumption), performance data could enable LLMs to rewrite communication materials to better resonate with their intended audience.

Our work only takes the first step in demonstrating the potential of RL augmented with aggregate performance feedback. We believe this is a promising and generalizable approach that bridges the gap between highly capable language models and tangible outcomes.

Acknowledgments

We gratefully acknowledge our close partnership on this project with the Monetization GenAI and Creative & Guidance teams at Meta: Yide Zhao, Yair Levi, Clare Zhang, Steven Barnett, Shenghong Wang, Meilei Jiang, Jerry Pan, Sanjian Chen, Shenxiu Liu, Zhonghua Qu, Xueting Yan, and Arghya Paul.

Author Contributions

A.N. and D.J. created reward model datasets and trained the reward model. D.J. and A.N. then post-trained the Text Generation LLM using RLPF, which resulted in the AdLlama model described in this paper. Y.C. curated the source datasets for both the reward model and Text Generation LLM. Y.C. and Y.B. managed the A/B testing process. D.J. drafted the initial version of the paper; all authors revised and reviewed the paper. Y.B. and Z.Z. advised on all aspects of the project.

Author Affiliations

All authors are either current or former employees of Meta. D.J., A.N., Y.C., and Y.B. are currently employed by Meta, while Z.Z. is a former employee. The work described in this paper by Z.Z. was conducted during Z.Z.'s employment at Meta.

References

- [1] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [2] OpenAI. Introducing ChatGPT, 2022. URL <https://openai.com/blog/chatgpt>.
- [3] Mina Lee, Percy Liang, and Qian Yang. CoAuthor: Designing a human-AI collaborative writing dataset for exploring language model capabilities. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19, 2022.
- [4] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [5] Eric Zhou and Dokyun Lee. Generative artificial intelligence, human creativity, and art. *PNAS Nexus*, 3(3):pgae052, 2024.
- [6] Harsh Kumar, David M Rothschild, Daniel G Goldstein, and Jake M Hofman. Math education with large language models: Peril or promise? *Available at SSRN 4641653*, 2023.
- [7] Hamsa Bastani, Osbert Bastani, Alp Sungu, Haosen Ge, Ozge Kabakcı, and Rei Mariman. Generative AI without guardrails can harm learning: Evidence from high school mathematics. *PNAS (forthcoming)*, 2025.
- [8] Pia Kreijkes, Viktor Kewenig, Martina Kuvalja, Mina Lee, Sylvia Vitello, Jake M Hofman, Abigail Sellen, Sean Rintel, Daniel G Goldstein, David M Rothschild, et al. Effects of LLM use and note-taking on reading comprehension and memory: A randomised experiment in secondary schools. *Available at SSRN 5095149*, 2025.
- [9] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez,

- Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature Medicine*, 29(8):1930–1940, 2023.
- [10] Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt, Narmin Ghaffari Laleh, Chiara Maria Lavinia Löffler, Sophie-Caroline Schwarzkopf, Michaela Unger, Gregory P Veldhuizen, et al. The future landscape of large language models in medicine. *Communications Medicine*, 3(1):141, 2023.
- [11] Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerová, et al. Adapted large language models can outperform medical experts in clinical text summarization. *Nature Medicine*, 30(4):1134–1142, 2024.
- [12] Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, et al. Medgemma technical report. *arXiv preprint arXiv:2507.05201*, 2025.
- [13] Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-level play in the game of Diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.
- [14] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [15] Google Deepmind AlphaProof and AlphaGeometry teams. AI achieves silver-medal standard solving international mathematical olympiad problems, 2024. URL <https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/>.
- [16] Erik Brynjolfsson. The Turing trap: The promise & peril of human-like artificial intelligence. *Daedalus*, 151(2):272–287, 2022.

- [17] Jenna Butler, Sonia Jaffe, Nancy Baym, Mary Czerwinski, Shamsi Iqbal, Kate Nowak, Sean Rintel, Abigail Sellen, Mihaela Vorvoreanu, Najeeb G. Abdulhazamid, Judith Amores, Reid Andersen, Kagonya Awori, Maxamed Axmed, danah boyd, James Brand, Georg Buscher, Dean Carignan, Martin Chan, Adam Coleman, Scott Counts, Madeleine Daepp, Adam Fournay, Daniel G. Goldstein, Andy Gordon, Aaron L Halfaker, Javier Hernandez, Jake Hofman, Jenny Lay-Flurrie, Vera Liao, Siân Lindley, Sathish Manivannan, Charlton Mcilwain, Subigya Nepal, Jennifer Neville, Stephanie Nyairo, Jacki O'Neill, Victor Poznanski, Gonzalo Ramos, Nagu Rangan, Lacey Rosedale, David Rothschild, Tara Safavi, Advait Sarkar, Ava Scott, Chirag Shah, Neha Parikh Shah, Teny Shapiro, Ryland Shaw, Auste Simkute, Jina Suh, Siddharth Suri, Ioana Tanase, Lev Tankelevitch, Adam Troy, Mengting Wan, Ryen W. White, Longqi Yang, Brent Hecht, and Jaime Teevan. Microsoft new future of work report 2023. *Microsoft Research Tech Report MSR- TR-2023-34*, 2023. URL <https://aka.ms/nfw2023>.
- [18] Manuel Hoffmann, Sam Boysel, Frank Nagle, Sida Peng, and Kevin Xu. Generative AI and the nature of work. Technical report, CESifo Working Paper, 2024.
- [19] Kunal Handa, Alex Tamkin, Miles McCain, Saffron Huang, Esin Durmus, Sarah Heck, Jared Mueller, Jerry Hong, Stuart Ritchie, Tim Belonax, et al. Which economic tasks are performed with AI? Evidence from millions of Claude conversations. *arXiv preprint arXiv:2503.04761*, 2025.
- [20] Harang Ju and Sinan Aral. Collaborating with AI agents: Field experiments on teamwork, productivity, and performance. *arXiv preprint arXiv:2503.18238*, 2025.
- [21] Erik Brynjolfsson, Danielle Li, and Lindsey Raymond. Generative AI at work. *The Quarterly Journal of Economics*, page qjae044, 2025.
- [22] Xiao Ni, Yiwei Wang, Tianjun Feng, Lauren Xiaoyuan Lu, Yitong Wang, and Congyi Zhou. Generative AI in action: Field experimental evidence on worker performance in e-commerce customer service operations. *Available at SSRN 5012601*, 2024.

- [23] Alexia Cambon, Brent Hecht, Ben Edelman, Donald Ngwe, Sonia Jaffe, Amy Heger, Mihaela Vorvoreanu, Sida Peng, Jake Hofman, Alex Farach, et al. Early LLM-based tools for enterprise information workers likely provide meaningful boosts to productivity. *Microsoft Research. MSR-TR-2023-43*, 2023.
- [24] Jonathan H Choi and Daniel Schwarcz. AI assistance in legal analysis: An empirical study. *Available at SSRN 4539836*, 2023.
- [25] Daniel Schwarcz, Sam Manning, Patrick Barry, David R Cleveland, JJ Prescott, and Beverly Rich. AI-powered lawyering: AI reasoning models, retrieval augmented generation, and the future of legal practice. *Available at SSRN 5162111*, 2025.
- [26] Fabrizio Dell’Acqua, Edward McFowland III, Ethan R Mollick, Hila Lifshitz-Assaf, Katherine Kellogg, Saran Rajendran, Lisa Kraye, François Cadelon, and Karim R Lakhani. Navigating the jagged technological frontier: Field experimental evidence of the effects of AI on knowledge worker productivity and quality. *Available at SSRN 4573321*, 2023.
- [27] Sida Peng, Eirini Kalliamvakou, Peter Cihon, and Mert Demirel. The impact of AI on developer productivity: Evidence from github copilot. *arXiv preprint arXiv:2302.06590*, 2023.
- [28] Zheyuan Kevin Cui, Mert Demirel, Sonia Jaffe, Leon Musolf, Sida Peng, and Tobias Salz. The effects of generative AI on high skilled work: Evidence from three field experiments with software developers. *Available at SSRN 4945566*, 2024.
- [29] Panagiotis Angelopoulos, Kevin Lee, and Sanjog Misra. Causal alignment: Augmenting language models with A/B tests. *Available at SSRN 4781850*, 2024.
- [30] Zenan Chen and Jason Chan. Large language model in creative work: The role of collaboration modality and user expertise. *Management Science*, 70(12):9101–9117, 2024.

- [31] Shakked Noy and Whitney Zhang. Experimental evidence on the productivity effects of generative artificial intelligence. *Science*, 381(6654):187–192, 2023.
- [32] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- [33] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [34] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.
- [35] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.
- [36] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- [37] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- [38] Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation

- framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- [39] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [40] Dentsu. Digital advertising expenditure worldwide from 2014 to 2027 (in billion U.S. dollars). Chart, Statista, December 2024. URL <https://www.statista.com/statistics/273717/global-internet-advertising-expenditure/>.
- [41] Dentsu. Share of digital in advertising revenue worldwide from 2023 to 2027. Graph, In *Statista*, December 2024. URL <https://www.statista.com/statistics/375008/share-digital-ad-spend-worldwide/>.
- [42] Meta. About Text Generation in Meta Ads Manager, 2024. URL <https://www.facebook.com/business/help/180641596861873>.
- [43] Meta. Create an ad with Text Generation in Meta Ads Manager, 2024. URL <https://www.facebook.com/business/help/497610041230617>.
- [44] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [45] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- [46] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

- [47] Meta. Creative best practices for text in ads, 2024. URL <https://www.facebook.com/business/help/223409425500940>.
- [48] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30, 2017.
- [49] Jayanta Mandi, Victor Bucarey, Maxime Mulamba Ke Tchomba, and Tias Guns. Decision-focused learning: Through the lens of learning to rank. In *International Conference on Machine Learning*, pages 14935–14947. PMLR, 2022.
- [50] Samuel Tan and Peter I Frazier. Asymptotically optimal regret for black-box predict-then-optimize. *arXiv preprint arXiv:2406.07866*, 2024.
- [51] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [52] Mark W Donoghoe and Ian C Marschner. Logbin: An R package for relative risk regression using the log-binomial model. *Journal of Statistical Software*, 86:1–22, 2018.
- [53] Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- [54] Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pages 10835–10866. PMLR, 2023.
- [55] Keyu Nie, Yinfei Kong, Ted Tao Yuan, and Pauline Berry Burke. Dealing with ratio metrics in A/B testing at the presence of intra-user correlation and segments. In *International Conference on Web Information Systems Engineering (WISE)*, pages 563–577, 2020.
- [56] Alex Deng, Ulf Knoblich, and Jiannan Lu. Applying the delta method in metric analytics: a practical guide with novel ideas. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 233–242, 2018.

[57] Edward L Frome. The analysis of rates using Poisson regression models. *Biometrics*, pages 665–674, 1983.

A Supplementary Materials

A.1 Training Data

A.1.1 Reward Model

We obtained our reward model (RM) training data from the Multiple Text Options feature in Meta Ads Manager [47]. This feature enables advertisers to manually submit various versions of body text, title text, and description text. Multiple Text Options then automatically tests and optimizes to find the ad variation that most effectively engages audiences. This feature depends entirely on *manually-written* text variations and therefore can be seen as a precursor to the AI-driven Text Generation feature (the focus of this paper).

Recall that the Text Generation feature targets the generation of *ad body text*. However, the CTR of an ad is not determined by its body text alone, but depends on multiple other factors, including the ad image, the ad’s targeting criteria, or the ad’s vertical (e.g., engagement on real estate ads is expected to be drastically different from restaurant ads). Thus, Multiple Text Options provided us with an important type of data: CTR data for ads that are *identical across all dimensions except for the ad body text*. This allows us to attribute the difference in CTR directly to the ad body text.

We filtered this data to include ad variations written in English with body text between 100 and 1000 characters (we excluded ads that were too short or too long) with at least 2,000 impressions on Facebook. From these data, we created preference pairs of ad variations where the two individual variations only differ in the ad body text, as previously illustrated in Figure 3. Our final RM training dataset included approximately 7 million preference pairs.

A.1.2 Language Model

The training data used for language model post-training (i.e., after the RM is trained) is sourced in the same way as described above, with the only exception being that preference pairs are not constructed. We refer to this as the “pointwise” version of the dataset (as opposed to “pairwise”). The pointwise dataset included approximately 5.5 million human-written text variation examples.

A.2 Reinforcement Learning with Performance Feedback

We adapt the RLHF paradigm by replacing pairwise human feedback with performance feedback. This change allows us to align the LLM not just with the preferences of a single human annotator (whose preferences might be noisy and not necessarily representative of the advertiser’s goals), but the ad’s CTR, a real-world performance metric that summarizes how the ad interacts with both the delivery system and Facebook users. Since the CTR is computed using all impressions of the ad, it becomes a less noisy estimate than a label from a handful of human annotators. The RLPF training pipeline has two main components: (1) reward model (RM) training and (2) post-training the LLM.

For RM training, we use the pairwise CTR preference dataset described above in Section A.1 of the supplementary materials, where preference pairs are decided by comparing CTRs over thousands of impressions (or more). We denote each row of data with a prompt x and a preference pair (y_w, y_l) . Following prior work on RLHF [48, 35], the RM is trained with an underlying Bradley-Terry preference assumption, namely that the probability of ad text y_1 being higher performing than ad text y_2 is given by

$$\begin{aligned}\mathbb{P}(y_1 \succ y_2 \mid x) &= \frac{\exp(r_\theta(x, y_1))}{\exp(r_\theta(x, y_1)) + \exp(r_\theta(x, y_2))} \\ &= \frac{1}{1 + \exp[-(r_\theta(x, y_1) - r_\theta(x, y_2))]} = \sigma(r_\theta(x, y_1) - r_\theta(x, y_2)),\end{aligned}$$

where $r_\theta(x, y)$ is the reward model with parameters θ and σ is the logistic function. The reward $r_\theta(x, y)$ represents the “strength” of response y given prompt x . Therefore, the Bradley-Terry model relates the preference to the relative strengths of two pieces of ad text. To fit the parameters θ , we

use maximum likelihood, arriving at the negative log-likelihood loss function:

$$\mathcal{L}_{\text{RM}}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}_{\text{RM}}} [\log \sigma(r_\theta(x, y_w) - r_\theta(x, y_l))],$$

where \mathcal{D}_{RM} is the pairwise preference dataset used for RM training. We use out-of-sample pairwise accuracy for hyperparameter tuning (learning rate, gradient accumulation), resulting in a model that reached approximately 57% out-of-sample pairwise accuracy.

After RM training, we apply the proximal policy optimization (PPO) algorithm [33, 51] to fine-tune the LLM and align it with high-performance ad text. A length penalty was added to the reward to counteract a tendency for a model to generate longer text, leading to the following PPO optimization formulation:

$$\max_{\pi_\phi} \mathbb{E}_{x \sim \mathcal{D}_{\text{LLM}}, y \sim \pi_\phi(y|x)} [r_\theta(x, y) - \beta \text{KL}[\pi_\phi(\cdot | x), \pi_{\text{ref}}(\cdot | x)] - \alpha \text{length}(y)],$$

where π_ϕ is the LLM, \mathcal{D}_{LLM} is the pointwise dataset used for LLM post-training, $\text{KL}(p, q)$ is the Kullback-Leibler divergence between two probability distributions p and q , $\text{length}(y)$ is the number of tokens in y , and β and α are weight parameters. We found that PPO can reliably increase the RM score, but subjective text quality started to decrease after a certain number of training steps (as evidenced by irrelevant or repetitive text). This is likely due to *overoptimization*, also known as *reward hacking*, a phenomenon where PPO starts to optimize the imperfections of the RM; see [54]. We deal with overoptimization by carefully selecting a model checkpoint using a combination of two strategies: (1) close monitoring the LLM training process using an *evaluation RM* trained on a different data split and (2) a small-scale human preference labeling.

A.3 Statistical Analysis

A.3.1 Descriptive Statistics and Covariates

In Table A.3.1, we present descriptive statistics of our sample of advertisers, along with balance of covariates. Below, we give details of each of the covariates used in the study.

Table A.3.1: Descriptive statistics of the advertisers in our A/B test. For categorical variables (is_business_account, budget_cat, expertise_cat, vertical), group differences were assessed using a χ^2 test. For the remaining variables, two-sample t-tests for equality of means were conducted, and two-sided p-values are reported. We do not observe statistically significant imbalances between the groups, indicating balanced sample characteristics.

Variable	Control <i>Imitation LLM v2</i>	Treatment <i>AdLlama</i>	p-value
<i>N</i>	17632	17217	
pre_exp_ctr (mean (SD))	0.027 (0.022)	0.027 (0.022)	0.766
pre_exp_engagement (mean (SD))	0.329 (2.556)	0.331 (4.411)	0.971
pre_exp_impressions (mean (SD))	16.497 (144.382)	15.811 (148.639)	0.662
pre_exp_ad_cnt (mean (SD))	0.448 (1.945)	0.454 (1.878)	0.753
account_age_yr (mean (SD))	3.238 (3.483)	3.276 (3.487)	0.307
nov_feb_ad_cnt (mean (SD))	2.007 (9.038)	2.114 (10.362)	0.304
nov_feb_variant_cnt (mean (SD))	5.983 (28.393)	6.435 (33.546)	0.175
is_business_account = 1 (%)	13468 (76.38)	13113 (76.16)	0.637
has_created_llm_ad = 1 (%)	4604 (26.11)	4570 (26.54)	0.366
is_new_advertiser = 1 (%)	1778 (10.08)	1729 (10.04)	0.912
budget_cat (%)			0.128
1.Low	8313 (47.15)	8031 (46.65)	
2.Mid	1901 (10.78)	1777 (10.32)	
3.High	7418 (42.07)	7409 (43.03)	
expertise_cat = 2.High (%)	2881 (16.34)	2710 (15.74)	0.131
vertical (%)			0.117
Advertising and Marketing	610 (3.46)	595 (3.46)	
Automotive	510 (2.89)	512 (2.97)	
Business to Business	353 (2.00)	437 (2.54)	
Consumer Packaged Goods	1195 (6.78)	1130 (6.56)	
Ecommerce	1681 (9.53)	1663 (9.66)	
Entertainment and Media	2330 (13.21)	2318 (13.46)	
Healthcare, Pharmaceuticals, and Biotech	778 (4.41)	716 (4.16)	
Other	1010 (5.73)	990 (5.75)	
Professional Services	2825 (16.02)	2782 (16.16)	
Publishing	635 (3.60)	551 (3.20)	
Restaurants	454 (2.57)	409 (2.38)	
Retail	3218 (18.25)	3176 (18.45)	
Technology	438 (2.48)	417 (2.42)	
Travel	581 (3.30)	570 (3.31)	
Unlisted	1014 (5.75)	951 (5.52)	

- `pre_exp_ctr`: Lifetime, pre-experiment CTR across all ads launched on Meta apps.
- `pre_exp_engagement`: Lifetime, pre-experiment engagement count across all ads launched on Meta apps. Units are in millions.
- `pre_exp_impressions`: Lifetime, pre-experiment impression count across all ads launched on Meta apps. Units are in millions.
- `pre_exp_ad_cnt`: Lifetime, pre-experiment ad count launched on all Meta apps. Units are in thousands.
- `account_age_yr`: The number of years since the advertiser’s account was created.
- `nov_feb_ad_cnt`: The number of ads created by the advertiser during the period after the launch of the initial Text Generation feature, but before the start of the A/B test.
- `nov_feb_variant_cnt`: The number of ad variants created by the advertiser during the period after the launch of the initial Text Generation feature, but before the start of the A/B test.
- `is_business_account`: A binary variable indicating whether the advertiser’s account is associated with a business page.
- `has_created_llm_ad`: A binary variable indicating whether the advertiser created an LLM-generated ad during the period after the launch of the initial Text Generation feature, but before the start of the A/B test.
- `is_new_advertiser`: A binary variable indicating whether the advertiser’s ads have combined for fewer than 1,000 total impressions on Meta apps (this includes many advertisers who have an account, but have delivered zero ad impressions).
- `budget_cat`: A categorical variable that groups advertisers based on their budget spend. The possible values are “1.Low,” “2.Mid,” and “3.High.”

- **expertise_cat**: A categorical variable that groups advertisers based on their expertise level with Meta’s ad system. High expertise means that the advertiser makes use of more advanced ad features. The possible values are “1.Low” and “2.High.”
- **vertical**: A categorical variable indicating the vertical in which the advertiser operates. Examples include “Retail,” “Travel,” or “Technology.”

A.3.2 Interpretation of Regression Specification under Log-link Function

Consider the log-binomial regression model specified in (1). Exponentiating both sides, we have

$$\begin{aligned}
\mathbb{E}(Y_i/n_i | X_i) &= \exp(\beta_2 \cdot \mathbf{1}_{\text{existing}(i)} \cdot \log(\text{CTR}_i^{\text{pre}}) + \beta_3 \cdot \mathbf{1}_{\text{new}(i)}) \exp(\text{“other covariates”}) \\
&= [\mathbf{1}_{\text{existing}(i)} (\text{CTR}_i^{\text{pre}})^{\beta_2} + \mathbf{1}_{\text{new}(i)} \exp(\beta_3)] \exp(\text{“other covariates”}) \\
&= (\text{“baseline” CTR}) \exp(\text{“other covariates”}).
\end{aligned}$$

The last equality illustrates that we can interpret the term in brackets as an estimated “baseline” CTR for advertiser i based on past performance. For existing advertisers (i.e., if $\text{existing}(i)$ is true), then baseline CTR term in brackets reduces to $(\text{CTR}_i^{\text{pre}})^{\beta_2}$, where we may interpret β_2 as an adjustment to convert between lifetime CTR on Meta’s apps to CTR on the Facebook platform. On the other hand, for new advertisers (i.e., $\text{new}(i)$ is true), we do not have an observation for their pre-experiment CTR and the term in brackets simply reduces to $\exp(\beta_3)$. This model-based term serves to estimate advertiser i ’s CTR in the absence of a pre-experiment CTR.

Therefore, we can interpret the entire regression as a “baseline CTR” further adjusted by other covariates. This interpretation is possible because we incorporated a log-transformed $\text{CTR}_i^{\text{pre}}$ into a regression formulation with a log link function.

A.3.3 Model-free Evidence

It is not immediately obvious how to compute the “average” CTR across the treatment and control groups—simply computing advertiser-level CTRs and then computing sample averages treats

low-impression advertisers and high-impression advertisers equivalently, but intuitively (absent a particular regression model), high-impression advertisers offer more precise CTR observations and should therefore be weighted higher.

It is common in A/B testing to compute the *global* CTR of each group by computing the ratio of the total number of engagements to the total number of impressions, i.e., $(\sum_i Y_i)/(\sum_i n_i)$. It turns out that global CTR is equivalent to an *impression-weighted* average, which matches our intuition above: $(\sum_i Y_i)/(\sum_i n_i) = \sum_i w_i (Y_i/n_i)$, where $w_i = n_i/(\sum_j n_j)$; see [55]. We apply the Delta method to compute the confidence interval around this estimate; see Equation 4 of [56] for an example. The model-free estimate is shown in Figure A.3.1.

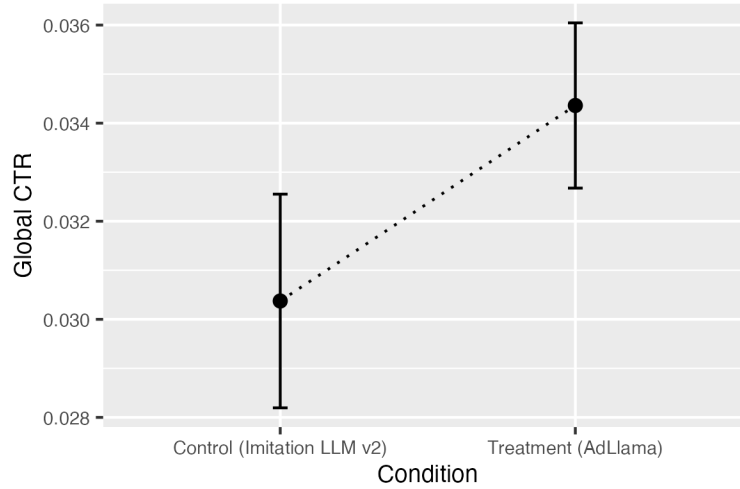


Figure A.3.1: **Model-free CTR estimates.** The control group’s CTR is estimated at 0.0304, with a 95% confidence interval of (0.0282, 0.0326). The treatment group’s CTR is 0.0344, with a 95% confidence interval of (0.0327, 0.0360), illustrated by the error bars in the plot. A two-sided z-test indicates a statistically significant difference between the groups ($p = 0.0046$).

A.3.4 Binomial, Logistic, and Poisson Regressions

First, the full results of our main log-binomial regression specification are given in Table A.3.2 (Table 1 is an abbreviated version). As an alternative to the log-binomial model of Equation 1, we also consider a logistic regression specification. The only difference from the log-binomial regression is a logit link function. These results are given in Table A.3.3; the results are in

agreement with the log-binomial specification.

So far, we have used heteroskedasticity-consistent standard errors in our models to ensure that inference is robust. Another approach is to use quasi-binomial regression, which introduces a dispersion parameter to adjust for variability exceeding that predicted by the binomial distribution, thereby providing more reliable standard error estimates. See Table A.3.4 for results. We observe that the standard errors under quasi-binomial specification are significantly less conservative than using robust standard errors. Overall, the results are qualitatively similar.

Finally, we consider and Poisson and quasi-Poisson regression specifications, which can be used to estimate *rates* using an “offset” term [57]. This is specified as follows:

$$\begin{aligned} \log \mathbb{E}(Y_i | X_i) = & \beta_0 + \beta_1 \cdot \text{AdLlama}_i + \beta_2 \cdot \mathbf{1}_{\text{existing}(i)} \cdot \log(\text{CTR}_i^{\text{pre}}) + \beta_3 \cdot \mathbf{1}_{\text{new}(i)} \\ & + \beta_{4:10}^T Z_i + \theta_{\text{vert}(i)} + \alpha_{\text{budget}(i)} + \lambda_{\text{expertise}(i)} + \log n_i. \end{aligned}$$

The term $\log n_i$ is called the “offset” and its coefficient is constrained to be 1, so that when moved to the left hand side, we obtain $\log(\mathbb{E}(Y_i | X_i)/n_i)$, precisely the CTR that we wish to estimate. The results for Poisson regression (with HC1 standard errors) are in Table A.3.5 and results for the quasi-likelihood variant are in Table A.3.6. The findings are again consistent with the other specifications.

A.3.5 Separate Engagement and Impression Linear Regressions

Alternatively, we can consider a regression specification where engagement and impressions are modeled separately as linear regressions, using the same covariates as we did above:

$$\begin{aligned} \text{Outcome}_i = & \beta_0 + \beta_1 \cdot \text{AdLlama}_i + \beta_2 \cdot \mathbf{1}_{\text{existing}(i)} \cdot \text{CTR}_i^{\text{pre}} + \beta_3 \cdot \mathbf{1}_{\text{new}(i)} \\ & + \beta_{4:10}^T Z_i + \theta_{\text{vert}(i)} + \alpha_{\text{budget}(i)} + \lambda_{\text{expertise}(i)} + \epsilon_i, \end{aligned}$$

where Outcome_i refers to either advertiser i ’s engagement count or impression count and ϵ_i is the error term. The results are given in Table A.3.7, where we find a statistically significant

Table A.3.2: The full version of Table 1 for the log-binomial regression; includes fixed effects.

	<i>Dependent variable:</i>
	engagement/impressions
treatment	0.0651** (0.0299)
log_pre_exp_ctr_existing	0.5931*** (0.0264)
is_new_advertiser	-2.3469*** (0.1569)
pre_exp_ad_cnt	0.0024 (0.0028)
pre_exp_impressions	-0.0001 (0.0001)
pre_exp_engagement	0.0036 (0.0049)
account_age_yr	0.0008 (0.0038)
nov_feb_ad_cnt	0.0013 (0.0019)
nov_feb_variant_cnt	-0.0003 (0.0006)
is_business_account	-0.0286 (0.0424)
has_created_llm_ad	-0.0109 (0.0345)
budget_cat: 2.Mid	0.0460 (0.0508)
budget_cat: 3.High	-0.0237 (0.0417)
expertise_cat: 2.High	0.0005 (0.0338)
vertical: Automotive	-0.1556 (0.1219)
vertical: Business to Business	-0.4097** (0.1867)
vertical: Consumer Packaged Goods	-0.1476 (0.1161)
vertical: Ecommerce	-0.0847 (0.1149)
vertical: Entertainment and Media	0.0648 (0.1134)
vertical: Healthcare, Pharmaceuticals, and Biotech	-0.1418 (0.1233)
vertical: Other	-0.2076 (0.1275)
vertical: Professional Services	-0.1612 (0.1248)
vertical: Publishing	-0.0107 (0.1451)
vertical: Restaurants	0.0117 (0.1185)
vertical: Retail	-0.1914 (0.1179)
vertical: Technology	-0.0369 (0.1360)
vertical: Travel	0.0163 (0.1179)
vertical: Unlisted	0.1067 (0.1681)
Constant	-1.1233*** (0.1506)
Observations	34,849

Note: HC1 robust standard errors in parentheses.

*p<0.1; **p<0.05; ***p<0.01

Table A.3.3: Logistic regression specification results. We report HC1 robust standard errors and two-sided p-values.

	<i>Dependent variable:</i>
	engagement/impressions
treatment	0.0685** (0.0311)
log_pre_exp_ctr_existing	0.6158*** (0.0277)
is_new_advertiser	-2.4418*** (0.1631)
pre_exp_ad_cnt	0.0024 (0.0029)
pre_exp_impressions	-0.0001 (0.0001)
pre_exp_engagement	0.0036 (0.0052)
account_age_yr	0.0009 (0.0040)
nov_feb_ad_cnt	0.0013 (0.0020)
nov_feb_variant_cnt	-0.0003 (0.0006)
is_business_account	-0.0318 (0.0444)
has_created_llm_ad	-0.0118 (0.0357)
budget_cat: 2.Mid	0.0484 (0.0528)
budget_cat: 3.High	-0.0271 (0.0433)
expertise_cat: 2.High	0.0006 (0.0351)
vertical: Automotive	-0.1572 (0.1275)
vertical: Business to Business	-0.4172** (0.1914)
vertical: Consumer Packaged Goods	-0.1497 (0.1213)
vertical: Ecommerce	-0.0848 (0.1202)
vertical: Entertainment and Media	0.0736 (0.1187)
vertical: Healthcare, Pharmaceuticals, and Biotech	-0.1415 (0.1286)
vertical: Other	-0.2132 (0.1330)
vertical: Professional Services	-0.1641 (0.1302)
vertical: Publishing	-0.0044 (0.1518)
vertical: Restaurants	0.0178 (0.1239)
vertical: Retail	-0.1952 (0.1232)
vertical: Technology	-0.0347 (0.1420)
vertical: Travel	0.0207 (0.1233)
vertical: Unlisted	0.1156 (0.1747)
Constant	-1.0005*** (0.1584)
Observations	34,849

Note: HC1 robust standard errors in parentheses.

*p<0.1; **p<0.05; ***p<0.01

Table A.3.4: Quasi-binomial regression specification results for both log and logit link functions. The fitted dispersion parameters are 8.427 and 8.433, respectively. We report two-sided p-values.

	<i>Dependent variable:</i>	
	engagement/impressions	
	<i>link = log</i>	<i>link = logit</i>
treatment	0.0651*** (0.0073)	0.0685*** (0.0076)
log_pre_exp_ctr_existing	0.5931*** (0.0069)	0.6158*** (0.0073)
is_new_advertiser	-2.3469*** (0.0339)	-2.4418*** (0.0354)
pre_exp_ad_cnt	0.0024* (0.0012)	0.0024* (0.0012)
pre_exp_impressions	-0.0001*** (0.00003)	-0.0001*** (0.00003)
pre_exp_engagement	0.0036*** (0.0011)	0.0036*** (0.0011)
account_age_yr	0.0008 (0.0010)	0.0009 (0.0011)
nov_feb_ad_cnt	0.0013* (0.0007)	0.0013* (0.0007)
nov_feb_variant_cnt	-0.0003 (0.0002)	-0.0003 (0.0002)
is_business_account	-0.0286** (0.0115)	-0.0318*** (0.0120)
has_created_llm_ad	-0.0109 (0.0085)	-0.0118 (0.0088)
budget_cat2.Mid	0.0460*** (0.0165)	0.0484*** (0.0172)
budget_cat3.High	-0.0237** (0.0105)	-0.0271** (0.0109)
expertise_cat2.High	0.0005 (0.0091)	0.0006 (0.0095)
vertical: Automotive	-0.1556*** (0.0298)	-0.1572*** (0.0310)
vertical: Business to Business	-0.4097*** (0.0316)	-0.4172*** (0.0327)
vertical: Consumer Packaged Goods	-0.1476*** (0.0257)	-0.1497*** (0.0268)
vertical: Ecommerce	-0.0847*** (0.0250)	-0.0848*** (0.0261)
vertical: Entertainment and Media	0.0648*** (0.0244)	0.0736*** (0.0255)
vertical: Healthcare, Pharmaceuticals, and Biotech	-0.1418*** (0.0323)	-0.1415*** (0.0337)
vertical: Other	-0.2076*** (0.0291)	-0.2132*** (0.0302)
vertical: Professional Services	-0.1612*** (0.0260)	-0.1641*** (0.0270)
vertical: Publishing	-0.0107 (0.0274)	-0.0044 (0.0286)
vertical: Restaurants	0.0117 (0.0333)	0.0178 (0.0349)
vertical: Retail	-0.1914*** (0.0248)	-0.1952*** (0.0258)
vertical: Technology	-0.0369 (0.0308)	-0.0347 (0.0319)
vertical: Travel	0.0163 (0.0278)	0.0207 (0.0290)
vertical: Unlisted	0.1067*** (0.0390)	0.1156*** (0.0405)
Constant	-1.1233*** (0.0355)	-1.0005*** (0.0373)
Observations	34,849	34,849

Note: Standard errors adjusted by the dispersion parameter.

*p<0.1; **p<0.05; ***p<0.01

Table A.3.5: Poisson regression results using log impressions as an offset; see Equation A.3.4. We report HC1 robust standard errors and two-sided p-values.

	<i>Dependent variable:</i>
	engagement/impressions
treatment	0.0658** (0.0299)
log_pre_exp_ctr_existing	0.5913*** (0.0263)
is_new_advertiser	-2.3431*** (0.1569)
pre_exp_ad_cnt	0.0023 (0.0028)
pre_exp_impressions	-0.0001 (0.0001)
pre_exp_engagement	0.0036 (0.0049)
account_age_yr	0.0008 (0.0038)
nov_feb_ad_cnt	0.0013 (0.0019)
nov_feb_variant_cnt	-0.0003 (0.0006)
is_business_account	-0.0299 (0.0425)
has_created_llm_ad	-0.0110 (0.0344)
budget_cat: 2.Mid	0.0461 (0.0506)
budget_cat: 3.High	-0.0259 (0.0417)
expertise_cat2.High	0.0001 (0.0338)
vertical: Automotive	-0.1515 (0.1224)
vertical: Business to Business	-0.4075** (0.1862)
vertical: Consumer Packaged Goods	-0.1442 (0.1165)
vertical: Ecommerce	-0.0812 (0.1153)
vertical: Entertainment and Media	0.0705 (0.1137)
vertical: Healthcare, Pharmaceuticals, and Biotech	-0.1364 (0.1233)
vertical: Other	-0.2064 (0.1277)
vertical: Professional Services	-0.1581 (0.1253)
vertical: Publishing	-0.0045 (0.1451)
vertical: Restaurants	0.0172 (0.1187)
vertical: Retail	-0.1868 (0.1182)
vertical: Technology	-0.0330 (0.1366)
vertical: Travel	0.0206 (0.1182)
vertical: Unlisted	0.1119 (0.1686)
Constant	-1.1313*** (0.1508)
Observations	34,849

Note: HC1 robust standard errors in parentheses.

*p<0.1; **p<0.05; ***p<0.01

Table A.3.6: Quasi-Poisson regression specification results. The fitted dispersion parameter is 8.127. We report two-sided p-values.

	<i>Dependent variable:</i>
	engagement/impressions
treatment	0.0658*** (0.0073)
log_pre_exp_ctr_existing	0.5913*** (0.0070)
is_new_advertiser	-2.3431*** (0.0340)
pre_exp_ad_cnt	0.0023* (0.0012)
pre_exp_impressions	-0.0001*** (0.00003)
pre_exp_engagement	0.0036*** (0.0011)
account_age_yr	0.0008 (0.0010)
nov_feb_ad_cnt	0.0013* (0.0007)
nov_feb_variant_cnt	-0.0003 (0.0002)
is_business_account	-0.0299*** (0.0115)
has_created_llm_ad	-0.0110 (0.0085)
budget_cat2.Mid	0.0461*** (0.0166)
budget_cat3.High	-0.0259** (0.0105)
expertise_cat2.High	0.0001 (0.0091)
vertical: Automotive	-0.1515*** (0.0299)
vertical: Business to Business	-0.4075*** (0.0316)
vertical: Consumer Packaged Goods	-0.1442*** (0.0258)
vertical: Ecommerce	-0.0812*** (0.0251)
vertical: Entertainment and Media	0.0705*** (0.0245)
vertical: Healthcare, Pharmaceuticals, and Biotech	-0.1364*** (0.0324)
vertical: Other	-0.2064*** (0.0291)
vertical: Professional Services	-0.1581*** (0.0261)
vertical: Publishing	-0.0045 (0.0275)
vertical: Restaurants	0.0172 (0.0335)
vertical: Retail	-0.1868*** (0.0249)
vertical: Technology	-0.0330 (0.0308)
vertical: Travel	0.0206 (0.0279)
vertical: Unlisted	0.1119*** (0.0391)
Constant	-1.1313*** (0.0357)
Observations	34,849

Note: Standard errors adjusted by the dispersion parameter. *p<0.1; **p<0.05; ***p<0.01

increase in engagement (+4.068, $p = 0.0023$) when using AdLlama, but no evidence for a change in impressions. Increased engagement while impressions are held constant is consistent with the result of increased CTR from the log-binomial regression of Table 1.

Table A.3.7: Separate linear regression results on engagements and impressions. All other notation remains consistent with Table 1. We report HC1 robust standard errors.

	<i>Dependent variable:</i>	
	engagement	impressions
treatment	4.068*** (1.334)	54.637 (49.730)
pre_exp_ctr_existing	175.495*** (31.954)	-1,982.590*** (692.366)
pre_exp_ad_cnt	1.248 (0.937)	44.755 (31.544)
pre_exp_engagement	0.407 (1.151)	-37.671 (31.945)
pre_exp_impressions	0.032* (0.019)	2.686*** (0.794)
account_age_yr	0.671*** (0.227)	21.168*** (8.094)
nov_feb_ad_cnt	0.255 (0.470)	-0.976 (13.499)
nov_feb_variant_cnt	0.081 (0.156)	4.943 (4.559)
is_new_advertiser	8.978*** (2.843)	174.461 (207.469)
is_business_account	5.757*** (1.143)	183.652*** (47.473)
has_created_llm_ad	5.677*** (1.847)	191.082*** (61.552)
budget_cat: 2.Mid	0.553 (1.557)	-32.907 (40.244)
budget_cat: 3.High	14.228*** (1.603)	396.337*** (54.707)
expertise_cat: 2.High	-1.183 (2.306)	21.213 (78.768)
vertical: Automotive	8.504** (4.031)	349.443** (162.753)
vertical: Business to Business	9.267 (7.137)	847.067 (548.643)
vertical: Consumer Packaged Goods	4.400 (4.069)	249.672 (153.272)
vertical: Ecommerce	7.378* (3.823)	265.690* (138.991)
vertical: Entertainment and Media	12.558*** (3.645)	295.972** (129.309)
vertical: Healthcare, Pharmaceuticals, and Biotech	1.378 (3.413)	57.893 (130.343)
vertical: Other	2.752 (3.158)	170.415 (124.527)
vertical: Professional Services	-0.467 (3.041)	43.760 (124.657)
vertical: Publishing	17.857*** (6.602)	410.544** (208.360)
vertical: Restaurants	3.119 (3.548)	75.335 (125.864)
vertical: Retail	4.415 (3.174)	197.073 (126.678)
vertical: Technology	9.659 (8.157)	376.426 (309.201)
vertical: Travel	14.820*** (4.668)	406.783** (166.148)
vertical: Unlisted	3.004 (3.623)	17.970 (190.378)
Constant	-11.150*** (3.526)	-140.400 (140.123)
Observations	34,849	34,849
R ²	0.017	0.015
Adjusted R ²	0.016	0.015

Note: HC1 robust standard errors in parentheses.

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

A.4 Imitation LLM Training Details

Imitation LLM v1 and v2 were developed internally at Meta prior to the RLPPF work described in this paper. Here, we give a brief description of the training process for these models.

- Both of these models used supervised fine-tuning (SFT), with the goal of imitating “good” ad text. This follows the instruction-tuning paradigm [33], where each row of training data is of the form (input ad text, target ad text). The input ad text represents the advertiser’s original ad text and the target ad text is a rewritten variation (i.e., an “improved” version). Running SFT on this data makes shifts the LLM’s response distribution to be more likely to generate the target ad text; in other words, “imitating” it.
- For Imitation LLM v1, based on 7B *Llama 2 Chat*, the target ad texts were synthetically distilled from the 70B *Llama 2 Chat* model, a more capable but larger model [44]. Instructions used to prompt the larger model are common ones used in ad copywriting, such as “paraphrase and shorten”, “make clear”, “make actionable”, “empathize”, “pose as a question”, or “focus selling point.”
- For Imitation LLM v2, in addition to synthetically distilled target ad texts, human-written ad texts were also added into the mix. This data was collected using similar instructions as for Imitation LLM v1. The synthetic dataset typically offers more creativity, but suffers from higher hallucination rates. On the other hand, the human-rewritten dataset oftentimes has higher quality, but may lack diversity. By training on both, Imitation LLM v2 struck a balance between creativity and quality.