

Paper Summary (Analytical)

Paper: Improving Generative Ad Text on Facebook using Reinforcement Learning
(arXiv:2507.21983v2, Dec 15, 2025)

1. Main Problem and Why It Matters

Large language models (LLMs) are powerful, but they often need post-training to perform well on real tasks. Most work on post-training focuses on subjective human preferences (e.g., RLHF). This paper asks a more concrete question: Can reinforcement learning (RL) with real performance metrics measurably improve a high-stakes, real-world product? The domain is online advertising, where small click-through rate (CTR) can represent large economic gains. The key contribution is to quantify the impact of RL post-training using actual advertiser outcomes rather than offline benchmarks.

2. Inputs and Outputs of the System

- Input: A human-written ad text provided by an advertiser.
- Output: Multiple AI-generated variations of the ad text, which the advertiser can review and edit before publishing.

The output is not just a single best rewrite. The system supports exploration by generating multiple options so the advertiser can select the most appropriate variant.

3. Data Used (Type, Source, Size)

The method relies on historical advertising performance data from Meta's "Multi-armed Bandit Optimization," where advertisers test multiple text variants while keeping other factors fixed (image, targeting, etc.). This makes it possible to attribute performance differences to the text itself.

Key properties:

- Type: Historical ad performance pairs (text variants with the same non-text attributes).
- Source: Meta/Facebook ad platform data.

- Size: In the live A/B test, the model was evaluated over ~35,000 advertisers and ~640,000 ad variations over 10 weeks.

4. Proposed Method in Simple Terms (RLPF)

The paper introduces Reinforcement Learning with Performance Feedback (RLPF).

1. Reward Model Training:

- Build pairs of ad texts where one performed better (higher CTR) than the other in multertext historical data.
- Train a reward model to predict preference: a text with higher CTR should get a higher reward score.

2. RL Fine-Tuning with PPO:

- Use the reward model as a proxy for real-world feedback.
- Fine-tune the LLM with Proximal Policy Optimization (PPO) to maximize reward while staying close to the base model.
- Add a length penalty to avoid overly long ad text.

This is conceptually between RLHF (human preference labels) and RL with verification (math/coding). Here the reward is a business metric (CTR).

Simple Pseudocode

5. Models Compared

- Imitation LLM v1: SFT on synthetic ad rewrites.
- Imitation LLM v2: SFT on synthetic + human-written rewrites.
- AdLlama (proposed): Imitation v2 + RLPF using reward model from performance feedback.

All models are based on Llama 2 Chat 7B.

6. Experiment Design (A/B Test)

- Period: Feb 16, 2024 – Apr 25, 2024 (10 weeks).
- Population: 34,849 US advertisers.
- Randomization: Advertiser-level assignment to control (Imitation v2) vs. treatment (AdLlama).
- Metrics: Advertiser-level CTR, total clicks, impressions, number of ads, number of ad variations.

The evaluation emphasizes advertiser-level ROI rather than per-ad metrics, consistent with business impact goals.

7. Main Results

- CTR improved by 6.7% relative ($p = 0.0296$) for AdLlama vs. Imitation v2.
- Absolute CTR: approximately 3.1% → 3.3%.
- Ad variations increased by 18.5%, suggesting higher adoption and satisfaction.
- Number of ads did not change, implying the effect is mainly in better text, not more ads.

These are meaningful gains in a mature ad platform where large improvements are hard to achieve.

8. Limitations

- Offline-only RL: The reward model is trained on historical data. It does not adapt to online or explore new trends.
- Single objective: Optimizes CTR only, not creativity, tone adherence, or brand constraints.
- Human selection not modeled: Advertisers choose which AI-generated variant to use; reward model does not include selection probability.

- Platform-level impacts: Broader effects (e.g., ad diversity, user experience) are analyzed.

9. Ideas for Continuation

- Online RL loop: incorporate real-time feedback from generated ads.
- Multi-objective reward: combine CTR with tone constraints, creativity, or advertiser satisfaction.
- Selection-aware reward: weight outcomes by likelihood of advertiser acceptance.
- Cross-domain adaptation: apply RLPF to email marketing, product description generation, customer support with performance metrics.

10. Why This Paper Is Important

This paper provides one of the first large-scale, real-world evaluations of RL policies for training impact on business metrics. It demonstrates that aligning LLMs with a mix of human and machine performance feedback (not just human preferences) can yield tangible economic value. This is a critical step in understanding how LLMs move from “capable” to “impactful.”