

Optimizing User Profiles via Contextual Bandits for Retrieval-Augmented LLM Personalization

Anonymous ACL submission

Abstract

Large Language Models (LLMs) excel at general-purpose tasks, but personalizing their responses to individual users remains challenging. Retrieval augmentation offers a lightweight alternative to fine-tuning by conditioning LLMs on user history records, yet existing strategies rely on heuristics (e.g., relevance to the query) that overlook the true contribution of records to personalization. To address these limitations, we propose PURPLE, a contextual bandit framework that optimizes User Profiles for LLM personalization. PURPLE operates as a re-ranking layer over candidate records, balancing efficiency with personalization quality. Across nine real-world personalization tasks spanning classification, regression, and short- and long-text generation, PURPLE consistently outperforms strong heuristic and retrieval-augmented baselines, establishing contextual bandit retrieval as a principled and scalable solution for personalized LLMs. Our code is available at: <https://anonymous.4open.science/r/ACL-2026-PURPLE-3096/>.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable success in various natural language processing tasks. As these models are increasingly applied to personalized applications, tailoring responses to individuals based on their own preferences has become a crucial challenge. Existing approaches for personalizing LLMs, such as Parameter-Efficient Fine-Tuning (PEFT) (Hu et al., 2022) and Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022), generally require modifying model parameters. These approaches incur high computational costs, demand frequent updates, and are impractical for real-time personalization at scale, especially when the LLM is not fully open-sourced. More-

over, continually fine-tuning models for different individuals would complicate safety evaluation and deployment, since each personalized variant would require separate testing.

In this paper, we study a lightweight approach to LLM personalization via retrieval augmentation (Wu et al., 2025), where user profiles are formed by retrieving and injecting past user records into the prompt. Prior work has shown that incorporating such profiles can effectively steer LLM outputs toward individual preferences (Salemi et al., 2024; Jiang et al., 2025). Compared to parameter-updating methods, retrieval-augmented personalization is lightweight, transparent, and readily deployable, as users can directly inspect and edit the records guiding generation. However, simply using the full user history risks introducing redundancy and noise and may exceed the model’s context window, while overly aggressive pruning may discard important personalization signals. This tension raises a central challenge: **which user records should be used to construct the user profile** so as to maximize *utility*, i.e., improve downstream task performance when conditioned on by the LLM?

Existing strategies for building user profiles often rely on heuristics, such as selecting records with the highest *relevance*, defined as semantic similarity to the query (Karpukhin et al., 2020). However, we argue that relevance is an unreliable proxy for personalization utility. This is because relevance-based heuristics conflate surface-level contextual similarity with the underlying preference signals that actually drive personalization. As a result, relevance scores emphasize how similar a record *looks* to the query not necessarily whether it supports the user’s current intent. To illustrate this dichotomy, consider Figure 1, where the user seeks “*relaxing movies for Friday night*”. Standard relevance heuristics prioritize record h_1 due to the high lexical overlap of “*Friday night*”, ig-

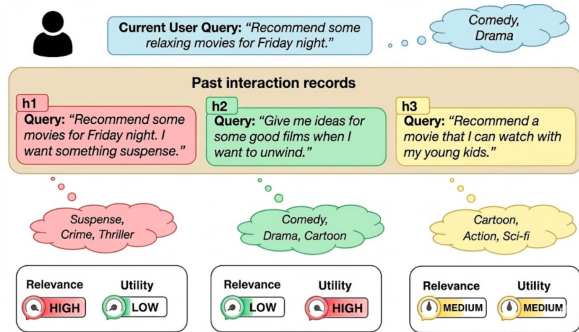


Figure 1: An illustration of the discrepancy between relevance and utility for a user’s past interaction records given a query. The thought bubbles illustrate the potential movie genres that satisfy the underlying user intent for each query.

noring the latent preference for “*suspense*” which contradicts the current need for relaxation. Conversely, record h_2 offers high utility by capturing the semantic intent to “*unwind*”, despite lacking surface-level keyword matches. This highlights that relevance scores prioritize how similar a record *looks* to the query, whereas utility depends on whether it supports the underlying user intent.

Beyond evaluating individual records, constructing an effective user profile requires set-aware reasoning rather than independent selection. Even if the utility of each record could be accurately estimated in isolation, a greedy top- k selection remains suboptimal. This is because record utilities do not compose additively; the contribution of a specific record depends heavily on the context of other retrieved records. Revisiting Figure 1, while h_2 and h_3 individually appear useful by reflecting a shared preference for lighthearted content, their combination is suboptimal. Jointly prompting with $\{h_2, h_3\}$ introduces conflicting constraints—mixing “*comedy*” preferences with “*young kids*” requirements—which dilutes generation quality compared to using h_2 alone. This demonstrates that personalization utility is non-monotonic: even valid preference signals can interfere with one another, meaning that simply accumulating “useful” records can actively degrade profile performance.

The limitations highlighted above call for a re-ranking mechanism that is directly optimized for downstream generation performance and is explicitly aware of interactions among records. Existing approaches fall short of these requirements: heuristic retrieval relies on static proxies like sim-

ilarity, while recent listwise rerankers, though capable of modeling dependencies, remain constrained by relevance-oriented supervision. To bridge this gap, we propose **PURPLE**, a framework that models user record selection as a contextual bandit problem (Langford and Zhang, 2007). In its formulation, the context consists of both the current query and the user’s past records. The selection policy is guided by a reward function reflecting downstream personalized text generation performance. PURPLE outputs a *propensity score* for each user record, which is passed through a Plackett-Luce ranking model to produce the final selected user records. This formulation enables the model to capture interactions between records and adaptively select those that are most beneficial. We train PURPLE end-to-end using the policy gradient method (Sutton et al., 1999).

Our main contributions are as follows:

- We introduce PURPLE, a framework that casts retrieval-augmented LLM personalization as a contextual bandit problem, adaptively optimizing user profiles beyond static heuristics.
- We show through extensive experiments on nine personalization tasks, covering classification, regression, and short- and long-text generation, that PURPLE consistently outperforms strong baselines in both effectiveness and efficiency.
- We perform comprehensive ablation studies on PURPLE’s key design choices and further validate PURPLE’s effectiveness through extended analysis and human evaluation.

2 Related work

LLMs for Personalization. LLMs demonstrate strong performance across domains (OpenAI, 2024), yet their outputs often diverge from user expectations because pre-training captures general rather than individual needs. Reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) and parameter-efficient finetuning (PEFT) (e.g., LoRA (Hu et al., 2022)) can align models with user preferences, but both require model finetuning and are impractical for end users who lack access to resources. A complementary direction personalizes LLMs through user profiles (Salemi et al., 2024), built from prior user interactions or external signals. Leveraging user profiles has shown benefits across multiple tasks requiring personalization, such as recommendation (Penaloza

et al., 2025), text summarization (Zhang et al., 2024), question answering (Wu et al., 2024), content generation (Shen et al., 2024), and personalized chatbot interaction (Jiang et al., 2025). Yet it remains unclear which user history records in a profile truly drive performance improvements, particularly in retrieval-augmented generation (RAG), where performance hinges on selecting semantically relevant context. Moreover, little analysis has been conducted on how to best select and compose user records into profiles with high personalization utility. Our work addresses this gap by studying how user profiles shape personalization in retrieval-augmented LLMs, and by proposing strategies for selecting user records to maximize downstream performance.

Retrieval-Augmented Language Models. Retrieval-augmented language models (RALMs) enhance parametric LMs with external memory to improve factuality and coverage. Early work such as REALM (Guu et al., 2020) and RAG (Lewis et al., 2020) jointly trained the retriever and LM, while Re2G (Glass et al., 2022) further incorporated a reranking module into this end-to-end pipeline. To reduce training costs, subsequent methods froze the LM and applied retrieval in-context. For example, In-Context RALM (Ram et al., 2023) leveraged LLMs for reranking, while REPLUG (Shi et al., 2024) distilled retrievers from LLMs. More recently, instruction-tuned variants such as SelfRAG (Asai et al., 2024) and RankRAG (Yu et al., 2024) jointly model retrieval and generation, but their reliance on large-scale finetuning renders them impractical for personalization.

The methods most relevant to our work are In-Context RALM and REPLUG; however, both incorporate only one retrieved record at a time, a limitation our method directly addresses. Specifically, REPLUG combines multiple records by weighting generation outputs with retrieval probabilities, while In-Context RALM periodically triggers retrieval during decoding at fixed steps and replaces previously used records. These designs arise because jointly reasoning over multiple records leads to a combinatorial explosion in the number of possible profiles. In contrast, our approach is explicitly designed to overcome this limitation by modeling cross-record dependencies and directly optimizing over multi-record profiles without resorting to such approximations.

LLMs for Reranking. Reranking methods can be categorized as pointwise, pairwise, or listwise. Pointwise models such as MonoBERT (Nogueira et al., 2019) and MonoT5 (Nogueira et al., 2020) score each query–document pair independently, while pairwise models such as DuoT5 (Pradeep et al., 2021) compare candidates in pairs. In contrast, listwise approaches jointly model the full candidate set and have recently been advanced by LLMs through prompt-only ranking (RankGPT (Sun et al., 2023)), model distillation (e.g., RankVicuna, RankZephyr, Lit5Distill, FIRST (Pradeep et al., 2023a,b; Tamber et al., 2023; Gangi Reddy et al., 2024)), and inference-time relevance extraction (ICR (Chen et al., 2025)). However, these methods conflate relevance with utility, which is not ideal for personalization. In this work, we instead train rerankers using downstream generation quality as feedback, prioritizing utility over semantic similarity.

3 Methodology

We formulate retrieval-augmented LLM personalization as a contextual bandit problem (Langford and Zhang, 2007), where the goal is to learn a policy that selects informative user records. Unlike classic multi-armed bandits, contextual bandits incorporate auxiliary information (e.g., the current query and user history) before making a selection, enabling direct optimization of retrieval strategies through policy gradient reinforcement learning. This aligns the selection of user records with downstream personalization objectives.

3.1 Problem Formulation

We consider a dataset $\mathcal{D} = \{(\mathcal{H}^u, x^u, y^u)\}_{u=1}^{|\mathcal{D}|}$, where each example consists of a user’s collection of history records \mathcal{H}^u , a query x^u , and a ground-truth personalized response y^u . Personalization is achieved by retrieving informative records from \mathcal{H}^u and supplying them as context to a frozen LLM, which then generates the final response. In practice, we apply PURPLE as a re-ranking module on top of a candidate pool selected by lightweight heuristics, ensuring low-latency inference compatible with large-scale systems. In the following development, we focus on a single user and omit the superscript index for brevity.

Let $\mathcal{H} = \{h_1, \dots, h_N\}$ denote the set of N history records for a user, where each record $h_i = (x_i, y_i)$ is an input–output pair (e.g., a query and

its answer from the user). Given a new query x , our goal is to construct a user profile from \mathcal{H} to condition the LLM for generating a personalized response. Formally, a profile is an ordered tuple $\mathcal{P} = \langle p_1, \dots, p_K \rangle \in \text{Perm}_K(\mathcal{H})$, which is a K -permutation of \mathcal{H} . We stress that the profile is order-sensitive: different permutations of the same K records correspond to distinct profiles and thus provide different inputs to the downstream LLM.

We formulate the selection of \mathcal{P} as a *contextual bandit problem*, where the context is given by the user’s history \mathcal{H} and the query x , and the action corresponds to selecting K records from \mathcal{H} to construct a profile. Formally, this formulation consists of the following key components:

- **Context:** $\mathcal{C} = (\mathcal{H}, x)$, where \mathcal{H} is the user’s collection of history records and x is the query. This representation captures both past user preferences and the immediate intent.
- **Actions:** $\mathcal{P} = \langle p_1, \dots, p_K \rangle \in \text{Perm}_K(\mathcal{H})$, which corresponds to selecting K distinct records from \mathcal{H} in a particular order. The action thus determines not only which records to use but also how they are arranged. The size of the action space is $N!/(N - K)!$.
- **Reward:** $R(\text{LLM}(\mathcal{P}, x), y)$, a function that measures the quality of the LLM-generated response $\text{LLM}(\mathcal{P}, x)$ relative to the ground-truth personalized response y .

We model the policy with a neural distribution $\pi_\theta(\cdot | \mathcal{C})$, parameterized by θ , which assigns probabilities to candidate user profiles given the context \mathcal{C} . The objective is to learn parameters θ such that the policy assigns higher probabilities to more informative profiles, which ultimately enhance personalized text generation. To this end, we maximize the expected reward over sampled user profiles, optimizing the following objective on a dataset \mathcal{D} spanning multiple users, each associated with a set of history records, a query, and the corresponding ground-truth answer:

$$\mathcal{J}(\theta) = \mathbb{E}_{(\mathcal{H}, x, y) \sim \mathcal{D}, \mathcal{P} \sim \pi_\theta(\cdot | \mathcal{C})} [R(\text{LLM}(\mathcal{P}, x), y)]. \quad (1)$$

It is challenging to directly optimize Equation 1 since the reward is not differentiable. To address this, we employ the likelihood ratio gradient estimator from reinforcement learning and stochastic

optimization (Williams, 1992; Sutton et al., 1999), which allows us to compute the gradient as:

$$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E}_{(\mathcal{H}, x, y) \sim \mathcal{D}, \mathcal{P} \sim \pi_\theta(\cdot | \mathcal{C})} [\nabla_\theta \log \pi_\theta(\mathcal{P} | \mathcal{C}) R(\text{LLM}(\mathcal{P}, x), y)]. \quad (2)$$

Since it is intractable to enumerate all profiles $\mathcal{P} \in \text{Perm}_K(\mathcal{H})$ during the optimization process, we estimate Equation 2 by randomly sampling $M = 32$ profiles. To stabilize training and reduce variance in gradient estimation, we apply z-score normalization over the rewards of these 32 profiles sampled for each example. The detailed gradient estimation scheme is provided in Equation 2 in Appendix A.

3.2 Model and Function Design

Design of $\pi_\theta(\cdot | \mathcal{C})$ Since different permutations of the selected records may lead to different final responses, we adopt the Plackett–Luce (PL) model, which assigns probabilities to profiles based on the scores of individual user records. Therefore, $\pi_\theta(\cdot | \mathcal{C})$ defines a distribution over all $N!/(N - K)!$ permutations of length K drawn from the user’s N history records. The probability assigned to a specific profile \mathcal{P} is given by:

$$\pi_\theta(\mathcal{P} | \mathcal{C}) = \prod_{k=1}^K \frac{f_\theta(p_k; \mathcal{C})}{S - \sum_{j=1}^{k-1} f_\theta(p_j; \mathcal{C})}, \quad (3)$$

where $S = \sum_{i=1}^N f_\theta(h_i; \mathcal{C})$, and $f_\theta(\cdot)$ is the user record encoder that outputs a propensity score in $[0, 1]$ for each record, indicating the model’s tendency to include that record in the user profile. During training, profiles are generated by sampling K records without replacement based on Equation 3. At inference time, the top- K records ranked by propensity scores are selected. Because our user record encoder is order-aware and rewards are assigned to ordered sets, the learned propensity score can be interpreted as each record’s contribution to the selected set. We further show in Sec. 5.3 that this ordering achieves higher final utility compared with other baselines.

Design of f_θ For the record encoder f_θ , we aim to capture the interdependencies among user records. A key design consideration is the trade-off between modeling dependencies at the token level versus the sentence level. While the former could, in principle, capture finer-grained interactions, it would quickly exceed the encoder’s context length. To address this, we adopt a late

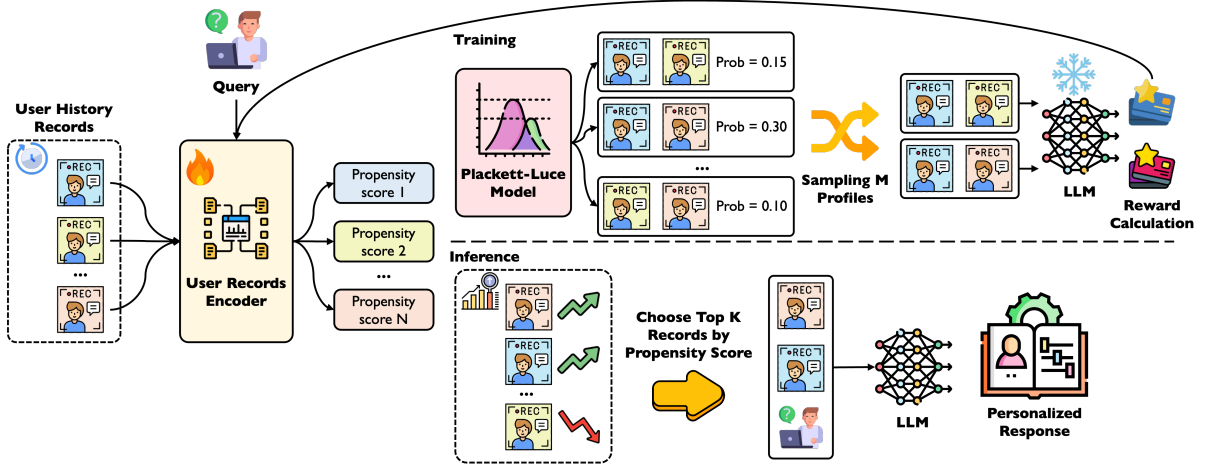


Figure 2: Workflow of the proposed PURPLE framework. User records encoder takes a user query and a list of user history records as input, outputting the propensity scores of all records. During **training**, a Plackett-Luce model is employed to convert the propensity scores to a probability distribution over all possible profiles, followed by sampling M profiles for gradient estimation. During **inference**, records with top K propensity scores are provided to the LLM along with the user query to generate a personalized response.

interaction strategy (Khattab and Zaharia, 2020), where we first obtain sentence-level embeddings with a pre-trained encoder, and then apply a Transformer encoder to model dependencies across records. Figure 2 illustrates the overall workflow of our method. Within the user record encoder, we utilize a pre-trained Contriever (Izacard et al., 2022) to obtain token embeddings for both the query and the records. Each record first cross-attends to the query at the token level, producing query-fused record embeddings that incorporate query information. A subsequent pooling operation is then applied to the updated record token embeddings to produce fixed-size sentence-level embeddings. embeddings are then processed by a Transformer encoder to model cross-record dependencies. We omit positional encodings to avoid ordering bias among records.

Design of Reward Function In this work, we propose an LLM-driven reward, where the policy is trained to maximize the log-likelihood that the LLM assigns to the target sequence. Formally, given a user profile \mathcal{P} , a query x , and a ground-truth personalized response y , we define the reward as:

$$\begin{aligned}
 R(\text{LLM}(\mathcal{P}, x), y) &= \log p_{\phi}(y \mid \mathcal{P}, x) \\
 &= \sum_{j=1}^{|y|} \log p_{\phi}(y_j \mid \mathcal{P}, x, y_{<j}),
 \end{aligned} \quad (4)$$

where ϕ are the parameters of the LLM and $p_{\phi}(\cdot)$ denotes its next-token distribution. Using the log-

likelihood of ground-truth sequences as the reward provides dense feedback signals, in contrast to downstream metrics such as accuracy, mean squared error, or ROUGE-1 (Liu et al., 2025). Moreover, we show in Appendix B that this objective is equivalent to maximizing the evidence lower bound (ELBO) of the marginalization-based RAG approach (Lewis et al., 2020), which becomes intractable in our setting due to the combinatorial explosion. In the next section, we empirically demonstrate that this log-likelihood-based reward is robust across diverse downstream tasks.

4 Experiments

4.1 Dataset and Evaluation

We evaluate the performance of PURPLE using Phi-4-Mini-Instruct (Microsoft, 2025) and Llama-3-8B-Instruct (Team, 2024) as the frozen LLM, and further scale to Llama-3-70B-Instruct (Team, 2024). Our experiments span a wide range of personalization settings, including personalized classification, regression, and both short- and long-text generation from the LaMP (Salemi et al., 2024) and LongLaMP (Kumar et al., 2024) benchmarks. We follow the prompt templates of Salemi et al. (2024) and Kumar et al. (2024) to incorporate user profiles into the original queries.

Specifically, we evaluate PURPLE on **nine personalization tasks**: two classification tasks — *Personalized Citation Identification* (Citation) and *Personalized Movie Tagging* (Movie), evalu-

Table 1: Results on the LaMP benchmark. The best and second-best results in each column are highlighted in **bold** and underlined, respectively. We report the mean and standard deviation over three runs with different random seeds for LLM generation. For GPT-5-nano, we report only a single run due to API cost constraints.

Task	Citation	Movie	Rating	News	Scholar	Tweet
Metric	Acc. / F1	Acc. / F1	MAE / RMSE	RG1 / RGL / MT	RG1 / RGL / MT	RG1 / RGL / MT
<i>With Phi-4-Mini-Instruct (3.84B)</i>						
BM25	63.9 _{0.53} / 63.8 _{0.69}	34.5 _{0.65} / 29.6 _{0.61}	0.438 _{0.00} / 0.852 _{0.01}	14.4 _{0.16} / 12.8 _{0.16} / 11.9 _{0.12}	39.7 _{0.08} / 33.2 _{0.04} / 42.2 _{0.08}	38.3 _{0.04} / 33.6 _{0.08} / 35.2 _{0.04}
Contriever	64.6 _{0.33} / 64.5 _{0.20}	36.0 _{0.37} / 31.1 _{0.24}	<u>0.424</u> _{0.01} / <u>0.830</u> _{0.03}	14.6 _{0.04} / 13.1 _{0.00} / 12.2 _{0.04}	39.7 _{0.00} / 33.3 _{0.08} / 42.0 _{0.08}	38.6 _{0.08} / 33.7 _{0.08} / 35.8 _{0.04}
IC-RALM	62.0 _{0.20} / 61.7 _{0.37}	33.6 _{0.12} / 28.8 _{0.16}	0.471 _{0.01} / 0.857 _{0.02}	13.4 _{0.04} / 11.9 _{0.04} / 11.0 _{0.00}	37.6 _{0.12} / 31.0 _{0.12} / 41.0 _{0.33}	38.1 _{0.12} / 33.3 _{0.16} / 35.2 _{0.12}
REPLUG	61.2 _{0.94} / 60.2 _{0.78}	<u>37.4</u> _{0.49} / <u>32.1</u> _{0.37}	0.486 _{0.01} / 0.891 _{0.02}	14.1 _{0.04} / 12.7 _{0.00} / 11.4 _{0.08}	38.1 _{0.63} / 32.8 _{0.90} / 40.4 _{1.51}	42.2 _{0.08} / 37.3 _{0.00} / 38.6 _{0.04}
RankGPT (Llama-3-8B-Instruct)	63.9 _{0.82} / 63.7 _{0.65}	34.2 _{0.94} / 28.9 _{1.18}	0.446 _{0.00} / 0.863 _{0.01}	14.7 _{0.33} / 13.1 _{0.24} / 12.3 _{0.24}	39.7 _{0.00} / 33.3 _{0.04} / 42.0 _{0.04}	38.2 _{0.04} / 33.4 _{0.02} / 35.2 _{0.08}
RankGPT (GPT-5-nano)	65.9 / 65.6	35.5 / 31.4	0.444 / 0.865	14.6 / 13.0 / 12.1	39.8 / 33.4 / 42.3	38.5 / 33.7 / 35.5
ICR (Llama-3-8B-Instruct)	65.2 _{0.53} / 65.0 _{0.45}	34.1 _{0.73} / 29.8 _{1.02}	<u>0.424</u> _{0.00} / <u>0.830</u> _{0.02}	<u>15.0</u> _{0.00} / <u>13.4</u> _{0.04} / <u>12.2</u> _{0.00}	39.5 _{0.08} / 33.1 _{0.08} / 42.0 _{0.04}	38.6 _{0.12} / 33.6 _{0.24} / 35.5 _{0.16}
PURPLE (Ours)	66.0 _{0.20} / 65.6 _{0.16}	38.6 _{0.33} / 34.2 _{0.53}	0.419 _{0.01} / 0.808 _{0.02}	15.1 _{0.04} / 13.5 _{0.00} / 12.6 _{0.08}	40.0 _{0.04} / 33.5 _{0.04} / 42.4 _{0.00}	<u>39.0</u> _{0.04} / <u>34.0</u> _{0.04} / <u>35.9</u> _{0.00}
<i>With Llama-3-8B-Instruct (8.03B)</i>						
BM25	56.9 _{0.61} / 56.5 _{0.53}	45.2 _{0.37} / 36.5 _{0.98}	0.328 _{0.01} / 0.664 _{0.02}	16.7 _{0.33} / 15.1 _{0.37} / 14.7 _{0.33}	41.0 _{0.02} / 35.2 _{0.04} / <u>40.6</u> _{0.02}	31.4 _{0.16} / 26.6 _{0.16} / 27.6 _{0.24}
Contriever	58.5 _{0.20} / <u>58.1</u> _{0.41}	47.2 _{0.37} / 39.1 _{0.29}	0.314 _{0.00} / 0.631 _{0.01}	<u>17.2</u> _{0.04} / <u>15.6</u> _{0.04} / <u>15.1</u> _{0.04}	41.1 _{0.08} / 35.4 _{0.12} / 40.5 _{0.04}	<u>32.1</u> _{0.16} / <u>27.2</u> _{0.24} / <u>28.5</u> _{0.16}
IC-RALM	<u>59.3</u> _{0.04} / 56.6 _{0.29}	38.0 _{0.82} / 29.9 _{0.41}	0.351 _{0.01} / 0.677 _{0.01}	14.1 _{0.24} / 12.5 _{0.24} / 12.1 _{0.08}	37.5 _{0.08} / 31.5 _{0.24} / 39.2 _{0.04}	29.2 _{0.69} / 24.6 _{0.57} / 25.6 _{0.49}
REPLUG	55.5 _{1.06} / 47.5 _{2.00}	42.0 _{1.35} / 31.8 _{1.14}	0.314 _{0.00} / 0.633 _{0.00}	14.8 _{0.04} / 13.2 _{0.00} / 11.8 _{0.04}	42.5 _{0.08} / 37.1 _{0.06} / 40.8 _{0.12}	30.4 _{0.20} / 26.1 _{0.12} / 25.8 _{0.37}
RankGPT (Llama-3-8B-Instruct)	57.0 _{0.20} / 56.3 _{0.12}	47.2 _{0.90} / 38.4 _{0.53}	0.318 _{0.01} / 0.637 _{0.01}	16.9 _{0.20} / 15.3 _{0.16} / 14.8 _{0.29}	41.0 _{0.04} / 35.4 _{0.08} / 40.5 _{0.16}	31.0 _{0.12} / 26.3 _{0.08} / 27.3 _{0.16}
RankGPT (GPT-5-nano)	59.5 / 58.0	45.1 / 36.2	0.321 / 0.638	17.1 / 15.4 / 15.0	41.0 / 35.3 / 40.5	31.5 / 26.5 / 27.8
ICR (Llama-3-8B-Instruct)	58.4 _{0.29} / 57.3 _{0.37}	<u>48.0</u> _{0.37} / <u>39.3</u> _{0.73}	<u>0.312</u> _{0.01} / <u>0.631</u> _{0.02}	17.1 _{0.04} / 15.4 _{0.00} / 14.9 _{0.00}	41.3 _{0.16} / 35.5 _{0.29} / 40.8 _{0.24}	31.8 _{0.33} / 26.8 _{0.24} / 28.1 _{0.29}
PURPLE (Ours)	59.2 _{0.82} / 58.8 _{0.78}	49.6 _{0.65} / 41.6 _{0.53}	0.307 _{0.01} / 0.624 _{0.01}	17.6 _{0.04} / 15.9 _{0.00} / 15.3 _{0.08}	<u>41.4</u> _{0.49} / <u>35.8</u> _{0.73} / 40.3 _{0.37}	32.5 _{0.04} / 27.5 _{0.00} / 28.8 _{0.04}
<i>With Llama-3-70B-Instruct (70.6B)</i>						
BM25	71.3 _{0.37} / 70.8 _{0.37}	54.5 _{0.37} / 47.2 _{0.41}	0.245 _{0.01} / 0.544 _{0.01}	18.1 _{0.33} / 16.5 _{0.29} / 14.9 _{0.37}	43.8 _{0.53} / 38.2 _{0.41} / 40.5 _{0.49}	36.3 _{0.16} / 30.9 _{0.20} / 33.0 _{0.12}
Contriever	71.2 _{0.78} / 70.8 _{0.73}	56.6 _{0.16} / 49.5 _{0.29}	<u>0.238</u> _{0.00} / <u>0.528</u> _{0.00}	18.6 _{0.08} / <u>17.0</u> _{0.08} / <u>15.6</u> _{0.08}	44.1 _{0.08} / <u>38.6</u> _{0.12} / 41.1 _{0.00}	36.5 _{0.00} / <u>31.4</u> _{0.00} / <u>33.3</u> _{0.00}
IC-RALM	66.5 _{0.04} / 66.4 _{0.00}	49.0 _{0.29} / 41.8 _{0.12}	0.261 _{0.00} / 0.555 _{0.00}	14.8 _{0.02} / 13.3 _{0.04} / 12.2 _{0.00}	39.7 _{0.00} / 34.1 _{0.24} / 38.6 _{0.16}	32.0 _{0.20} / 27.4 _{0.00} / 28.8 _{0.20}
REPLUG	65.8 _{0.33} / 65.4 _{0.41}	51.7 _{0.04} / 44.0 _{0.08}	0.258 _{0.01} / 0.554 _{0.02}	15.2 _{0.04} / 13.9 _{0.04} / 12.1 _{0.04}	41.8 _{0.12} / 36.6 _{0.33} / 38.6 _{0.35}	32.1 _{0.08} / 27.6 _{0.12} / 27.7 _{0.08}
RankGPT (Llama-3-8B-Instruct)	69.7 _{0.16} / 69.1 _{0.12}	<u>56.9</u> _{0.16} / 48.8 _{0.41}	0.247 _{0.00} / 0.546 _{0.01}	18.0 _{0.24} / 16.5 _{0.29} / 15.2 _{0.29}	44.0 _{0.04} / <u>38.6</u> _{0.08} / 41.1 _{0.12}	35.8 _{0.00} / 30.6 _{0.00} / 32.4 _{0.08}
RankGPT (GPT-5-nano)	73.8 / 73.5	55.3 / 48.2	0.240 / <u>0.523</u>	18.7 / 17.0 / 15.8	44.6 / 38.8 / 41.5	<u>36.6</u> / 31.3 / 33.2
ICR (Llama-3-8B-Instruct)	71.6 _{0.16} / 71.0 _{0.16}	56.8 _{0.24} / <u>49.6</u> _{0.57}	<u>0.238</u> _{0.00} / 0.531 _{0.00}	18.3 _{0.00} / 16.8 _{0.04} / 15.3 _{0.16}	44.3 _{0.06} / 38.5 _{0.04} / <u>41.2</u> _{0.20}	36.3 _{0.16} / 31.1 _{0.24} / 33.2 _{0.20}
PURPLE (Ours)	<u>72.2</u> _{0.49} / <u>71.8</u> _{0.57}	57.0 _{0.08} / 49.9 _{0.41}	0.236 _{0.02} / 0.520 _{0.01}	18.8 _{0.04} / 17.1 _{0.00} / <u>15.6</u> _{0.08}	<u>44.4</u> _{0.02} / 38.8 _{0.06} / 41.0 _{0.04}	37.3 _{0.01} / 32.1 _{0.04} / 34.0 _{0.06}

ated with Accuracy and F1; one regression task — *Personalized Product Rating* (Rating), evaluated with MAE and RMSE; and six generation tasks, evaluated with ROUGE-1 (RG1), ROUGE-L (RGL) (Lin, 2004), and METEOR (MT) (Banerjee and Lavie, 2005). The generation tasks are further divided into short-text generation — *Personalized News Headline Generation* (News), *Personalized Scholarly Title Generation* (Scholar), and *Personalized Tweet Paraphrasing* (Tweet) — and long-text generation — *Personalized Abstract Generation* (Abstract), *Personalized Topic Generation* (Topic), and *Personalized Product Review Generation* (Review). In all experiments, we first use Contriever (Izcard et al., 2022) to retrieve the top 20 records as the user history \mathcal{H} , and then select 5 of them to construct the user profile \mathcal{P} . Appendix C further elaborates implementation details.

4.2 Baseline Methods

We focus on the setting where the LLM is kept frozen and no ground-truth profile is available for training the reranker. Therefore, we compare with three categories of prior methods that, likewise, neither fine-tune the LLM nor rely on supervision from ground-truth retrieval results.

The baselines we compare with include **(i) Zero-Shot Rerankers**, represented by

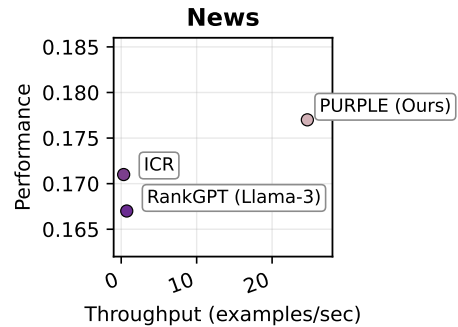


Figure 3: Performance-throughput graph on the News dataset. PURPLE is faster than LLM-based rerankers while achieving better performance.

RankGPT (Sun et al., 2023) and ICR (Chen et al., 2025). For both methods, we adopt Llama-3-8B-Instruct as the reranker LLM. We also report the performance of RankGPT with GPT-5-nano to reflect an upper bound of the methods that distill from the ranking results of state-of-the-art proprietary LLMs (Pradeep et al., 2023a,b; Tamber et al., 2023; Gangi Reddy et al., 2024). **(ii) In-Context Retrieval-Augmented Language Models**, represented by IC-RALM (Ram et al., 2023) and REPLUG (Shi et al., 2024). Both methods consider only one record at a time when generating a response. They incorporate multiple records from the user profile either through marginalization

Table 2: Ablation studies of PURPLE using Phi-4-Mini-Instruct. CA and RDM stand for cross-attention and record dependency modeling, respectively. MR refers to using task-specific evaluation metric as the reward.

Task	Citation	Movie	Rating	News	Scholar	Tweet
Metric	Acc. / F1	Acc. / F1	MAE / RMSE	RG1 / RGL / MT	RG1 / RGL / MT	RG1 / RGL / MT
PURPLE	66.2 / 65.8	38.2 / 33.6	0.405 / 0.788	15.2 / 13.5 / 12.5	40.0 / 33.5 / 42.4	39.1 / 34.0 / 35.9
w/o CA	64.8 / 64.5	35.1 / 29.7	0.440 / 0.816	14.8 / 13.2 / 12.4	40.0 / 33.5 / 42.2	39.1 / 34.1 / 36.0
w/o RDM	61.3 / 60.6	35.0 / 31.1	0.449 / 0.850	14.5 / 12.8 / 11.9	39.7 / 33.1 / 41.9	39.0 / 34.0 / 36.1
w/ MR	64.8 / 64.8	38.0 / 33.3	0.433 / 0.854	15.0 / 13.2 / 12.4	39.5 / 32.9 / 41.8	38.7 / 33.7 / 35.6

(REPLUG) or context switching (In-Context RALM). Additionally, we include (iii) **Efficient Sparse and Dense Retrievers**, applied directly as rerankers. We include BM25 (Robertson and Zaragoza, 2009) for the sparse retriever and Contriever (Izacard et al., 2022) for the dense retriever. These methods represent the efficiency-oriented side of the efficiency–performance trade-off.

5 Experiment Results

5.1 Overall Performance Comparison

Table 1 presents the results of PURPLE and baseline methods on the LaMP benchmark, while Table 3 presents the results on the LongLaMP benchmark. The main findings are as follows:

PURPLE consistently outperforms strong baselines across LLM scales Across all tasks and LLMs of varying sizes, PURPLE achieves consistent improvements over existing methods. Compared with Contriever, which is of comparable model size, our learned propensity scores provide more effective ranking signals than raw relevance. Compared with zero-shot rerankers, which use much larger backbone LLMs and incur higher inference cost, PURPLE achieves stronger personalization with a much smaller model, since reinforcement learning training allows us to better capture the utility of profiles formed by multiple records. Compared with in-context RALMs, which provide user records one at a time to the LLM and combine multiple records post hoc, our single-stage modeling more effectively captures personalized signals, highlighting the advantage of treating user profiles holistically.

PURPLE outperforms strong baselines with high computational throughput. Figure 3 shows that, on a representative LaMP dataset, PURPLE outperforms existing methods while being more efficient. PURPLE achieves higher

performance than ICR and RankGPT (Llama-3) while maintaining high computational throughput.

PURPLE is effective across task types, including regression. As our reward is based on the log-probability that the LLM assigns to the ground-truth response, it does not directly reflect numerical distances between regression targets. Nevertheless, PURPLE still achieves strong gains on the regression task (Rating). This demonstrates that log-probability provides a principled and broadly applicable reward signal across diverse task formats.

5.2 Ablation Studies

Table 2 presents ablation studies of PURPLE using Phi-4-Mini-Instruct. Overall, we examine three key design choices. First, instead of performing token-level cross attention, we test a simplified variant that encodes the entire query into a single embedding and prepends it as an extra token to the Transformer encoder (w/o CA). This approach is less effective, indicating that fine-grained token-level interactions between the query and user records are crucial for accurate personalization. Second, we remove the Transformer encoder entirely, resulting in a point-wise scoring model where each record is processed independently (w/o RDM). This variant shows the largest performance drop across tasks. Although the model can still leverage individually informative records, it fails to model dependencies such as redundancy and complementarity among records. Third, we examine an alternative reward design that uses task-specific evaluation metrics as the reward (Accuracy, MAE, and ROUGE-1). This metric-based reward is substantially more sparse and coarse-grained, providing weaker learning signals and making optimization more challenging. Empirically, this variant consistently underperforms PURPLE across all tasks, thereby

538 strongly motivating our reward design choice.

539 These results highlight that token-level cross at-
540 tention, cross-record dependency modeling, and
541 the log-probability-based reward are all indispens-
542 able, validating our design choice of treating user
543 profiles as structured contexts rather than isolated
544 records, and providing denser reward signals via
545 ground-truth log-probabilities given by the LLM.

5.3 Analysis: Selecting topK at inference

547 To further examine the quality of the learned
548 propensity scores, we compare the top-5 selec-
549 tions of PURPLE, ICR, RankGPT, and Contriever
550 on Citation, Rating, News, and Tweet. For each
551 test example, we consider the top-5 records pro-
552 posed by each method and randomly sample five
553 orderings from the $5! = 120$ permutations as con-
554 trols. As shown in Figure 4, orderings induced
555 by our learned propensity scores are more fre-
556 quently ranked as the best among the six candi-
557 dates. This result indicates that our scoring func-
558 tion better captures relative dependencies between
559 records, rather than relying on local pairwise re-
560levance alone. These findings highlight that our
561 method not only identifies useful records but also
562 arranges them in an order that maximizes down-
563 stream personalization utility.

5.4 Human Evaluation

565 To complement our automatic evaluation and ver-
566 ify the practical utility of the selected profiles,
567 we conduct a rigorous side-by-side human eval-
568 uation. While automatic metrics measure lexical
569 overlap, they often fail to capture the subtle stylis-
570 tic nuances, such as voice, stance, and informal
571 phrasing, that define true personalization. We fo-
572 cus our human evaluation on *Personalized Tweet*
573 *Paraphrasing* (Tweet), which requires significant
574 stylistic adaptation to user personas. We randomly
575 sample 250 instances. For each instance, human
576 evaluators are presented with the input query, a
577 snippet of the user’s history, and two anonymized
578 model outputs: one generated using one of the
579 most recent baseline ICR and one using PURPLE.
580 Evaluators perform a blind, forced-choice compar-
581 ison to determine which output better reflected the
582 user’s specific persona while maintaining seman-
583 tic fidelity. The results demonstrate that PURPLE
584 achieves a decisive victory, surpassing ICR by a
585 margin of +14.4% (57.2% vs. 42.8%). This con-
586 firms that our method effectively captures personal
587 stylistic features, such as slang usage and infor-

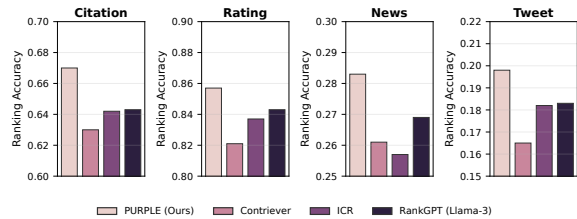


Figure 4: Ranking accuracy across various tasks. PURPLE achieves the highest accuracy on the examined datasets, consistently outperforming heuristic retrievers, LLM-based rerankers, and in-context rerankers.

588 mal tone, where heuristic relevance retrieval often
589 fails. We provide detailed qualitative case studies
590 analyzing these behaviors in Appendix D.

6 Conclusion

592 In this work, we study the problem of retrieval-
593 augmented personalization for large language
594 models. Our key intuitions are twofold: (i)
595 *record relevance does not reliably predict person-*
596 *alization utility*, and (ii) *utility is non-monotonic*
597 across records, making greedy aggregation sub-
598 optimal. To address these limitations, we pro-
599 pose **PURPLE**, a contextual bandit framework
600 that optimizes user profiles by directly leverag-
601 ing downstream performance as feedback. PUR-
602 PLE jointly models query–record interactions and
603 cross-record dependencies, enabling adaptive se-
604 lection of user profiles beyond static heuristics. In-
605 stead of raw task-specific metrics, PURPLE relies
606 on LLM-based probability as the reward, provid-
607 ing denser reward signals that facilitate more ef-
608 fective training. Extensive experiments on nine
609 real-world personalization tasks spanning clas-
610 sification, regression, and text generation show
611 that PURPLE consistently outperforms heuristic
612 retrievers, LLM-based rerankers, and in-context
613 RALMs, while being significantly more efficient.
614 These results establish contextual bandit-based re-
615 trieval as a principled and scalable paradigm for
616 personalizing LLMs.

617 Limitations

618 While PURPLE demonstrates strong effectiveness
619 across diverse personalization tasks, it has several
620 limitations. First, the method relies on the log-
621 likelihood of ground-truth personalized responses
622 under a frozen LLM as the training reward, which
623 assumes access to high-quality personalized su-
624 pervision. In practical deployment, such explicit
625 supervision may be sparse or unavailable, with
626 user feedback often taking implicit or noisy forms
627 (e.g., engagement signals), and extending PUR-
628 PLE to these settings remains an open challenge.
629 Second, although using log probability as the re-
630 ward provides a unified and task-agnostic opti-
631 mization signal across classification, regression,
632 and generation tasks, PURPLE is still trained sep-
633 arately for each dataset and task in our current
634 study. We do not explicitly evaluate the extent to
635 which a policy learned under this unified objec-
636 tive can generalize across tasks or domains, and
637 exploring more transferable training and deploy-
638 ment strategies is left for future work.

References 639

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *The Twelfth International Conference on Learning Representations*. 640
641
642
643
644
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics. 645
646
647
648
649
650
651
652
- Shijie Chen, Bernal Jimenez Gutierrez, and Yu Su. 2025. Attention in Large Language Models Yields Efficient Zero-Shot Re-Rankers. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net. 653
654
655
656
657
658
- Revanth Gangi Reddy, JaeHyeok Doo, Yifei Xu, Md Arafat Sultan, Deevya Swain, Avirup Sil, and Heng Ji. 2024. **FIRST: Faster improved listwise reranking with single token decoding**. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8642–8652, Miami, Florida, USA. Association for Computational Linguistics. 659
660
661
662
663
664
665
666
- Michael Glass, Gaetano Rossiello, Md Faisal Mahub Chowdhury, Ankita Naik, Pengshan Cai, and Alfio Gliozzo. 2022. **Re2G: Retrieve, Rerank, Generate**. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2701–2715, Seattle, United States. Association for Computational Linguistics. 667
668
669
670
671
672
673
674
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pappas, and Mingwei Chang. 2020. Retrieval Augmented Language Model Pre-Training. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3929–3938. PMLR. 675
676
677
678
679
680
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *International Conference on Learning Representations*. 681
682
683
684
- Edward J. Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*. 685
686
687
688
689
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. *Transactions on Machine Learning Research*. 690
691
692
693
694

695	Bowen Jiang, Zhuoqun Hao, Young-Min Cho, Bryan Li, Yuan Yuan, Sihao Chen, Lyle Ungar, Camillo J. Taylor, and Dan Roth. 2025. Know me, respond to me: Benchmarking llms for dynamic user profiling and personalized responses at scale . <i>Preprint</i> , arXiv:2504.14225.	752
696		753
697		754
698		755
699		
700		
701	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781, Online. Association for Computational Linguistics.	756
702		757
703		758
704		
705		
706		
707		
708		
709	Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and effective passage search via contextualized late interaction over BERT . In <i>Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , Sigir '20, pages 39–48, New York, NY, USA. Association for Computing Machinery.	759
710		760
711		761
712		762
713		763
714		764
715		
716	Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization . <i>Preprint</i> , arXiv:1412.6980.	765
717		766
718		767
719	Ishita Kumar, Snigdha Viswanathan, Sushrita Yerra, Alireza Salemi, Ryan A. Rossi, Franck Dernoncourt, Hanieh Deilamsalehy, Xiang Chen, Ruiyi Zhang, Shubham Agarwal, Nedim Lipka, Chien Van Nguyen, Thien Huu Nguyen, and Hamed Zamani. 2024. LongLaMP: A Benchmark for Personalized Long-form Text Generation .	768
720		769
721		
722		
723		
724		
725		
726	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with PagedAttention . In <i>Proceedings of the 29th Symposium on Operating Systems Principles</i> , Sosp '23, pages 611–626, New York, NY, USA. Association for Computing Machinery.	770
727		771
728		772
729		773
730		774
731		775
732		776
733		777
734		778
735	John Langford and Tong Zhang. 2007. The Epoch-Greedy algorithm for contextual multi-armed bandits. In <i>Proceedings of the 21st International Conference on Neural Information Processing Systems</i> , NIPS'07, pages 817–824, Red Hook, NY, USA. Curran Associates Inc.	779
736		780
737		781
738		782
739		783
740	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 9459–9474. Curran Associates, Inc.	784
741		785
742		786
743		787
744		788
745		789
746		790
747		791
748	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	792
749		793
750		794
751		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805

806	Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024. LaMP: When Large Language Models Meet Personalization . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7370–7392, Bangkok, Thailand. Association for Computational Linguistics.	861
807		862
808		863
809		864
810		865
811		866
812		867
813	Xiaoteng Shen, Rui Zhang, Xiaoyan Zhao, Jieming Zhu, and Xi Xiao. 2024. PMG: Personalized multimodal generation with large language models . In <i>Proceedings of the ACM Web Conference 2024, Www '24</i> , pages 3833–3843, New York, NY, USA. Association for Computing Machinery.	868
814		869
815		870
816		871
817		872
818		873
819	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. REPLUG: Retrieval-Augmented Black-Box Language Models . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8371–8384, Mexico City, Mexico. Association for Computational Linguistics.	874
820		875
821		876
822		877
823		878
824		879
825		880
826		881
827		882
828		883
829	Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agents . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 14918–14937, Singapore. Association for Computational Linguistics.	884
830		885
831		886
832		887
833		888
834		889
835		890
836		891
837	Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. In <i>Advances in Neural Information Processing Systems</i> , volume 12. MIT Press.	892
838		893
839		894
840		895
841		896
842	Manveer Singh Tamber, Ronak Pradeep, and Jimmy Lin. 2023. Scaling down, LiTting up: Efficient zero-shot listwise reranking with seq2seq encoder-decoder models . <i>Preprint</i> , arXiv:2312.16098.	897
843		898
844		899
845		900
846	Llama Team. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	901
847		902
848	Ronald J. Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning . <i>Machine Learning</i> , 8(3):229–256.	903
849		904
850		905
851	Bin Wu, Zhengyan Shi, Hossein A. Rahmani, Varsha Ramineni, and Emine Yilmaz. 2024. Understanding the role of user profile in the personalization of large language models . <i>Preprint</i> , arXiv:2406.17803.	906
852		907
853		908
854		909
855	Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can Chen, Ye Yuan, Lianming Huang, Xue Liu, Tei-Wei Kuo, Nan Guan, and Chun Jason Xue. 2025. Retrieval-augmented generation for natural language processing: A survey . <i>Preprint</i> , arXiv:2407.13193.	910
856		911
857		912
858		913
859		914
860		915
	Yue Yu, Wei Ping, Zihan Liu, Boxin Wang, Jiaxuan You, Chao Zhang, Mohammad Shoeybi, and Bryan Catanzaro. 2024. RankRAG: Unifying Context Ranking with Retrieval-Augmented Generation in LLMs . In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 121156–121184. Curran Associates, Inc.	916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960

A Details of Gradient Estimation

To estimate the gradient in Equation 2, we first draw a batch of examples $\{(\mathcal{H}_b, x_b, y_b)\}_{b=1}^B$. For each example, we sample M profiles $\mathcal{P}_b^1, \dots, \mathcal{P}_b^M$ from $\pi_\theta(\cdot | \mathcal{C}_b)$, and finally compute the empirical mean. This learning procedure corresponds to the REINFORCE algorithm (Sutton et al., 1999), with gradient estimate:

$$\nabla_\theta \mathcal{J}(\theta) \approx \frac{1}{B} \sum_{b=1}^B \frac{1}{M} \sum_{m=1}^M \nabla_\theta \log \pi_\theta(\mathcal{P}_b^m | \mathcal{C}_b) \tilde{r}_b^m. \quad (5)$$

To reduce variance in gradient estimation, we apply reward normalization over the M profiles sampled for each example. Concretely, for each example with rewards $\mathbf{r}_b = [r_b^1, \dots, r_b^M]^\top$, where $r_b^m = R(\Phi(\mathcal{P}_b^m, x_b), y_b)$, the normalized reward is computed as $\tilde{r}_b^m = \frac{r_b^m - \text{mean}(\mathbf{r}_b)}{\text{std}(\mathbf{r}_b)}$.

B Motivating our Reward

The specific choice of using the log-probability of ground truth personalized response is grounded in the generative modeling perspective of retrieval-augmented generation (RAG) (Lewis et al., 2020), where the user profile is treated as a latent variable and the response likelihood is obtained by marginalizing over all possible profile selections. Applying Jensen’s inequality to the training objective in this setting gives:

$$\begin{aligned} & \mathbb{E}_{(\mathcal{H}, x, y) \sim \mathcal{D}} \\ & \left[\log \sum_{\mathcal{P} \in \text{Perm}_K(\mathcal{H})} \pi_\theta(\mathcal{P} | \mathcal{C}) p_\Phi(y | \mathcal{P}, x) \right] \quad (6) \\ & \geq \mathbb{E}_{(\mathcal{H}, x, y) \sim \mathcal{D}, \mathcal{P} \sim \pi_\theta(\cdot | \mathcal{C})} [\log p_\Phi(y | \mathcal{P}, x)]. \end{aligned}$$

Therefore, maximizing the expected reward under our reinforcement learning objective is equivalent to maximizing the evidence lower bound (ELBO), with p_Φ modeled by a frozen LLM.

C Implementation Details

We employ a frozen pre-trained Contriever to first encode both queries and user history records into token embeddings. The only trainable components are the remaining modules of the user records encoder. These include a cross-attention layer that integrates query information into record embeddings, a Transformer encoder that captures

inter-record dependencies, and an MLP decoder that maps the updated record encodings into scalar propensity scores. We set the number of Transformer encoder layers to $l = 12$, resulting in a parameter size roughly twice that of Contriever, while still being substantially smaller than the zero-shot rerankers and in-context RALMs used as baselines. For gradient estimation, we use a batch size of $B = 16$ and sample $M = 32$ user profiles for each example. We train the model for 10 epochs using the Adam optimizer (Kingma and Ba, 2017) with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a learning rate of 1×10^{-4} . During training, we apply a gradient clipping norm of 1.0. The checkpoint achieving the best validation performance is selected for testing.

In all experiments, we use frozen LLMs both to generate personalized responses and to evaluate the log-probability of ground-truth responses conditioned on the query and user profiles (i.e., our reward). For generation, we set the temperature to $T = 0.7$ and employ nucleus sampling (Holtzman et al., 2020) with $top_p = 0.8$. For Phi-4-Mini-Instruct and Llama-3-8B-Instruct, we deploy on a single NVIDIA H100 GPU. For Llama-3-70B-Instruct, we deploy the model across four NVIDIA H100 GPUs using vLLM (Kwon et al., 2023). All LLMs are deployed in BF16 precision. Training of PURPLE is conducted on the same GPUs used for LLM deployment.

D Qualitative Case Studies

In this section, we present detailed case studies from our human evaluation to illustrate the qualitative superiority of PURPLE over the baseline ICR. We present one representative example where the two methods diverge significantly.

In the Tweet task, preserving the logic of the original sentence is as critical as matching the user’s informal voice. The example below demonstrates a failure in logical reasoning by the baseline. The input tweet contains a comparative structure (“at least we attended... unlike someone”). The baseline fails to parse this dependency, resulting in a paraphrase that states the exact opposite of the truth. PURPLE successfully disentangles the complex sentence structure and preserves the correct meaning.

Table 3: Results on the LongLaMP benchmark. The best and second-best results in each column are highlighted in **bold** and underlined, respectively. We report the mean and standard deviation over three runs with different random seeds for LLM generation. For GPT-5-nano, we report only a single run due to API cost constraints.

Task	Abstract	Topic	Review
Metric	R1 / RL / MT	R1 / RL / MT	R1 / RL / MT
<i>With Phi-4-Mini-Instruct (3.84B)</i>			
BM25	38.6 _{0.12} / 22.0 _{0.16} / 26.2 _{0.08}	25.2 _{0.41} / 12.7 _{0.20} / <u>17.1</u> _{0.16}	25.8 _{1.39} / 13.2 _{0.53} / 16.5 _{0.16}
Contriever	38.5 _{0.04} / 21.9 _{0.12} / 26.1 _{0.08}	<u>25.7</u> _{1.80} / <u>13.0</u> _{0.73} / 16.6 _{0.20}	25.2 _{1.92} / 13.1 _{0.57} / 15.9 _{0.78}
IC-RALM	37.2 _{0.04} / 21.1 _{0.02} / 25.1 _{0.04}	24.8 _{1.47} / 12.4 _{0.78} / 15.9 _{0.20}	25.0 _{1.35} / 12.7 _{0.49} / 15.8 _{0.16}
REPLUG	36.0 _{0.24} / 21.2 _{0.20} / 23.8 _{0.20}	23.5 _{1.06} / 12.1 _{0.20} / 13.1 _{0.41}	20.8 _{2.94} / 11.2 _{1.14} / 12.5 _{1.63}
RankGPT (Llama-3-8B-Instruct)	38.6 _{0.12} / 22.0 _{0.08} / 26.2 _{0.08}	25.1 _{0.53} / 12.7 _{0.29} / <u>17.1</u> _{0.18}	25.1 _{1.63} / 12.9 _{0.53} / 15.8 _{0.45}
RankGPT (GPT-5-nano)	39.1 / 22.4 / 26.9	24.9 / 12.5 / 17.5	<u>27.1 / 13.7 / 16.6</u>
ICR (Llama-3-8B-Instruct)	38.5 _{0.24} / 22.0 _{0.16} / 26.2 _{0.16}	<u>25.7</u> _{1.71} / <u>13.0</u> _{0.73} / 16.5 _{0.24}	25.8 _{1.67} / 13.2 _{0.53} / 16.4 _{0.45}
PURPLE (Ours)	<u>38.8</u> _{0.08} / <u>22.1</u> _{0.16} / <u>26.4</u> _{0.12}	26.2 _{1.14} / 13.2 _{0.61} / 17.5 _{0.12}	27.9 _{0.04} / 14.4 _{0.08} / 17.1 _{0.04}
<i>With Llama-3-8B-Instruct (8.03B)</i>			
BM25	42.2 _{0.04} / <u>24.2</u> _{0.04} / 31.8 _{0.08}	29.8 _{0.73} / 14.6 _{0.24} / <u>20.2</u> _{0.12}	<u>30.9</u> _{2.00} / <u>15.3</u> _{0.82} / <u>20.8</u> _{0.41}
Contriever	<u>42.3</u> _{0.24} / <u>24.2</u> _{0.29} / 31.7 _{0.24}	<u>30.6</u> _{1.39} / 15.2 _{0.53} / 20.0 _{0.32}	30.1 _{2.49} / 14.9 _{1.02} / 19.7 _{0.90}
IC-RALM	38.3 _{0.90} / 20.4 _{0.78} / 29.1 _{0.37}	25.0 _{0.90} / 11.4 _{1.06} / 15.5 _{1.92}	29.8 _{1.27} / 13.6 _{0.98} / 18.3 _{0.98}
REPLUG	38.2 _{0.45} / 20.7 _{0.33} / 28.4 _{0.20}	20.2 _{1.22} / 10.8 _{0.57} / 11.9 _{1.27}	15.5 _{2.04} / 8.8 _{0.94} / 8.8 _{1.02}
RankGPT (Llama-3-8B-Instruct)	42.2 _{0.04} / <u>24.2</u> _{0.04} / 31.9 _{0.04}	30.8 _{1.39} / 15.2 _{0.57} / 20.4 _{0.36}	30.5 _{2.45} / 15.1 _{1.02} / 20.4 _{0.82}
RankGPT (GPT-5-nano)	42.5 / 24.5 / 32.1	28.7 / 14.2 / <u>20.2</u>	33.6 / 16.5 / 21.4
ICR (Llama-3-8B-Instruct)	41.9 _{0.24} / 23.9 _{0.20} / 31.5 _{0.16}	<u>30.6</u> _{1.35} / <u>15.1</u> _{0.57} / <u>20.2</u> _{0.20}	30.7 _{1.96} / 15.2 _{0.78} / 20.3 _{0.41}
PURPLE (Ours)	<u>42.3</u> _{0.08} / <u>24.2</u> _{0.12} / <u>32.0</u> _{0.24}	<u>30.4</u> _{1.63} / <u>15.0</u> _{0.73} / <u>19.8</u> _{0.24}	<u>30.2</u> _{1.61} / <u>14.9</u> _{0.83} / <u>20.0</u> _{0.94}

User Persona / History:

- @tommogger thanks tom, and glad it all went well
- @callamKAPOW I tried to change it but the pic was too big and I can't find the old one
- Thinks I can pass English, I just wish I could be as confident about geography

Input Tweet:

Paraphrase the following tweet without any explanation before or after it: @picklesticks6 and I am in the same boat, but at least we attended the math lessons unlike someone we both know who is taking their exam today and feels sad about it. =(

Ground Truth Target:

@picklesticks6 same, but at least we turned up to the maths lessons unlike someone we know whos taking their exam today =(

ICR (Baseline) Generation:

@picklesticks6 we both missed math classes, and it's making our exam prep tougher, unlike someone who went and feels down about it today. =(

PURPLE (Ours) Generation:

@picklesticks6 feels the same about math exams as me, but at least we went to the lessons, unlike someone else who feels sad today about their exam. =(

Analysis:

This case highlights a critical semantic failure in the baseline. **ICR** generates "we both missed math classes," which directly contradicts the input stating "at least we attended." It effectively swaps the roles of the speaker and the third party. **PURPLE** demonstrates superior natural language understanding by correctly preserving the fact that the speaker *did* attend the lesson. Furthermore, **PURPLE** better matches the user's persona (which discusses academic confidence in English/Geography) by maintaining the conversational flow without losing the factual premise.

E LLM Usage

During the preparation of this manuscript, we employed large language models (LLMs) to assist with English writing refinement and style polishing. All technical content, including the design of **PURPLE**, theoretical formulations, experimental setup, and reported results, was conceived, implemented, and validated by the authors. The LLMs were used solely for linguistic improvement and did not contribute to the research methodology or experimental findings.