

# A Framework for Personalized Persuasiveness Prediction via Context-Aware User Profiling

Anonymous ACL submission

## Abstract

Estimating the persuasiveness of messages is critical in various applications, from recommender systems to safety assessment of LLMs. While it is imperative to consider the target persuadee’s characteristics, such as their values, experiences, and reasoning styles, there is currently no established systematic framework to optimize leveraging a persuadee’s past activities (e.g., conversations) to the benefit of a persuasiveness prediction model. To address this problem, we propose a context-aware user profiling framework with two trainable components: a *query generator* that generates optimal queries to retrieve persuasion-relevant records from a user’s history, and a *profiler* that summarizes these records into a profile to effectively inform the persuasiveness prediction model. Our evaluation on the ChangeMyView Reddit dataset shows consistent improvements over existing methods across multiple predictor models, with gains of up to +13.77%p in F1 score. Further analysis shows that effective user profiles are context-dependent and predictor-specific, rather than relying on static attributes or surface-level similarity. Together, these results highlight the importance of task-oriented, context-dependent user profiling for personalized persuasiveness prediction.<sup>1</sup>

## 1 Introduction

Large language models (LLMs) are increasingly used in decision-support applications that aim to influence human behavior or beliefs, such as health coaching, tutoring, and targeted marketing (Salvi et al., 2024a; Hackenburg et al., 2025a). In these settings, an LLM may generate or evaluate multiple candidate messages (e.g., campaign messages for marketing companies) to assist a human decision maker, requiring the system to determine which message is most likely to persuade a target user.

<sup>1</sup>We will release our code and data upon publication.

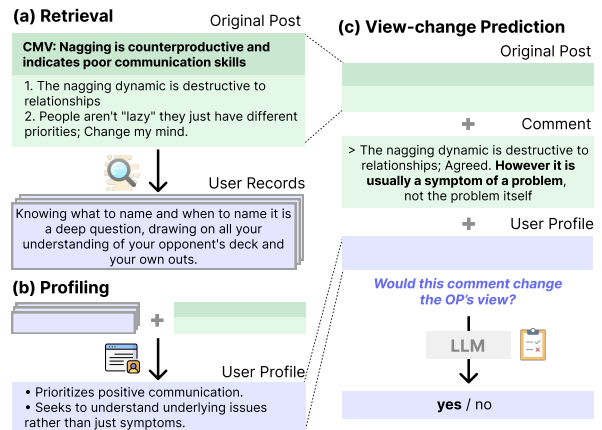


Figure 1: Overview of the view-change prediction task with context-aware user profiling on CMV. Given an original post, the system (a) retrieves relevant user records, (b) constructs a textual user profile, and (c) predicts whether a comment will change the user’s view.

We refer to this problem as *persuasiveness prediction*, defined as predicting a user’s belief or attitude change in response to a given message (Perloff, 2021). The main challenge in persuasiveness prediction stems from the fact that persuasion is inherently personalized: the same argument may be compelling for one user but ineffective for another, depending on factors such as beliefs, values, experiences, and reasoning style (Lukin et al., 2017a; Durmus and Cardie, 2018a; Khatib et al., 2020). As a result, accurate persuasiveness prediction requires inferring how each individual user is likely to interpret and respond to a message. In practice, such inference must rely on signals from a user’s past interaction history, as explicit user attributes are often unavailable. This motivates methods that infer user characteristics from historical interactions to enable personalized persuasiveness prediction.

We formulate this problem as a *personalized view-change prediction* task using data from the ChangeMyView (CMV) Reddit forum. In this setting, each instance consists of (i) an original post

expressing a user’s view on a specific topic, (ii) a comment responding to the post with the intent to change that view, and (iii) the user’s historical *records*, consisting of past Reddit posts and comments. The goal is to predict whether the comment will successfully change the user’s stated view by leveraging user information inferred from the user’s history (Figure 1c; Tan et al., 2016). To make accurate predictions, a system requires to construct user representations from their historical records that capture what matters for the current persuasion context (Figure 1a,b; Li et al., 2016; Xu et al., 2025). However, existing approaches typically rely on heuristic retrieval methods (e.g., selecting recent records or random sampling) to select relevant historical records and generic profiling techniques (e.g., extracting demographic traits) to summarize user characteristics from those records (Hackenburg et al., 2025b; Al Khatib et al., 2020; Salvi et al., 2024b). We argue that these approaches are insufficient because persuasion is inherently context-dependent: which aspects of a user’s history are informative depends on the topic and stance of the original post as well as the argument in the candidate message (Tan et al., 2016; Ji et al., 2018).

To address this limitation, we propose a framework with two learnable modules for generating user profile in textual form (Figure 1) : (i) a *query generator* that produces retrieval queries to identify persuasion-relevant records from a user’s history, and (ii) a *profiler* that summarizes the retrieved records into a textual user profile, conditioned on the original post. This profile, together with the original post, is then used by a predictor model to determine whether the candidate message would change the user’s view. Examples of the final user profiles are presented in Appendix E.

We train the components in the following order to provide clear learning signals: 1) First, we train the *profiler* using Direct Preference Optimization (DPO) (Rafailov et al., 2023): given random combinations of user records and the original post, the profiler learns to generate textual profiles that capture user characteristics useful for view-change prediction in that persuasion context. 2) Second, we derive record-level persuasion utility scores—measures of how useful each individual user record is for predicting view change in a given context—which are used to train the query generator. Because direct supervision on the usefulness of individual records is unavailable, we estimate each

record’s usefulness by evaluating the task performance across multiple random combinations of a user’s records; records that frequently appear in high-performing combinations receive higher scores. These scores serve as a learning signal that encourages the query generator to retrieve records most valuable for persuasiveness prediction (Ghorbani and Zou, 2019). 3) Finally, we train the query generator via DPO by treating queries that retrieve high-utility records as *chosen* and those that fail to do so as *rejected*, optimizing the model to generate queries that maximize persuasion utility.

Our evaluation on the CMV dataset shows consistent gains over prior approaches, demonstrating the effectiveness of our framework for personalized persuasiveness prediction. Further analyses show that persuasion-relevant user characteristics vary across posts and predictor models, highlighting the need for predictor-specific, context-aware user profiles rather than generic or static attributes.

Our contributions are threefold:

1. We propose a framework that trains retrieval and user profiling modules without ground-truth annotations, using view-change prediction performance as the supervision signal.
2. We introduce a *persuasion-aware query generation* method that enables retrieval of user records relevant for personalized persuasion.
3. We provide empirical evidence that effective user profiling is context-dependent and predictor-specific.

## 2 Related Work

**Early Work on Personalized Persuasion** While Tan et al. (2016) established the view-change prediction task on the CMV dataset, subsequent work has enriched modeling by incorporating richer linguistic features, such as interaction dynamics and discourse relations (Ji et al., 2018; Hidey and McKeeown, 2018). Another line of research demonstrates the importance of personalization by leveraging persuadee characteristics in persuasion outcome prediction, including ideology, demographic, and personality traits (Lukin et al., 2017b; Durmus and Cardie, 2018b, 2019a,b; Al Khatib et al., 2020). However, they largely rely on pre-defined, explicit user attributes. Leveraging recent advances in LLMs, we infer richer persuasion-relevant user information from the users’ past writings.

**User Profiling for LLM Personalization** Early work has formulated LLM personalization as making models behave *like* a specific user given their historical writings (Salemi et al., 2024b; Mysore et al., 2024). Building on this, several studies have explored retrieval- and profiling-based approaches (Richardson et al., 2023; Li et al., 2024; Salemi et al., 2024a; Zhang, 2024). These methods focus on linguistic style and topical relevance, remaining limited in capturing user *values* (Qin et al., 2025), which are crucial for personalized persuasion. Studies on personalized dialogue agents construct user profiles via summarization to generate user-aligned responses (Zhong et al., 2024; Wang et al., 2025), but remain limited in dynamically adapting profiles to the current interaction context.

Another line of work fine-tunes the predictor on users’ historical data to encode user-specific information (Zhang et al., 2025). While effective, it requires retraining as new user data arrives, limiting scalability in practice. In contrast, we keep the predictor fixed and focus on constructing context-aware user profiles, and thus do not directly compare against such approaches.

### 3 User Profiling Framework

#### 3.1 Problem Formulation

We construct a dataset from the ChangeMyView (CMV) Reddit forum, where users post opinions and award a *delta* to comments that change their views. The discussions cover diverse topics, including politics, personal values, and everyday issues. To support personalized prediction, we collect each user’s historical posts and comments from both CMV and other subreddits. We split this dataset into training, validation, and test sets with an 8:1:1 ratio. For training and validation, we subsample up to 100 user records per post. This is necessary because user records, which are scored via repeated view-change prediction (Section 3.3.2), are large and highly variable in size (mean 784, max 19K records per post). Concretely, we use the *delta* comment as a retrieval query and build the pool using a hybrid retriever that combines BM25 with BGE-M3 semantic similarity. This substantially reduces computational cost while preserving records that are informative for view-change prediction.

To formalize the task, we represent each data instance as a tuple  $(u, x_i, c_i, y_i, R_u)$ . Here,  $u$  denotes a user;  $x_i$  is an original post authored by  $u$  in the CMV forum that expresses the user’s ini-

tial view on a topic; and  $c_i$  is a comment written by another user in response to  $x_i$ . Each comment is labeled as a *delta* or *non-delta*, with a label  $y_i \in \{0, 1\}$  indicating whether it changed the user’s view ( $y_i = 1$  if a *delta* was awarded). We further provide access to the user’s historical records  $R_u = \{r^{u,1}, r^{u,2}, \dots, r^{u,|R_u|}\}$ , where each  $r^{u,j}$  is a Reddit post or comment written by  $u$  prior to  $x_i$ . The personalized view-change task is formulated as  $\tilde{y}_i = f(u, R_u, x_i, c_i)$ , where the goal is to predict whether  $c_i$  will change user  $u$ ’s view expressed in  $x_i$ , given the user’s history  $R_u$ .

#### 3.2 Inference

To address this task, we introduce a three-stage inference pipeline comprising *retrieval*, *profiling*, and *view-change prediction* (Figure 2a–c). Since  $R_u$  is typically large and noisy, direct conditioning on the entire set is impractical. Instead, we construct a compact user profile  $P_i$  that summarizes persuasion-relevant information about  $u$  and condition the prediction solely on  $P_i$ . The construction of  $P_i$  consists of two stages: retrieval and profiling.

The retrieval stage (Figure 2a) selects a subset of  $k$  records from  $R_u$  that are directly used for profile construction. We first generate a retrieval query  $q_i$  using a trainable query generator  $\phi^{\text{query}}$ , which takes the original post  $x_i$  as input. Using an embedding-based retriever  $\mathcal{M}^{\text{ret}}$ , we retrieve the top- $k$  records most relevant to  $q_i$ :

$$\{r^{u,i_1}, \dots, r^{u,i_k}\} = \mathcal{M}^{\text{ret}}(q_i, R_u, k) \subseteq R_u.$$

The profiling stage (Figure 2b) constructs a natural-language user profile  $P_i$  by summarizing the retrieved records into a textual representation that the predictor model can effectively utilize. We employ a trainable LLM-based profiler  $\phi^{\text{prof}}$  that takes the retrieved records and the original post  $x_i$  as input, enabling the profile to be conditioned on the persuasion context expressed in  $x_i$ :

$$P_i = \phi^{\text{prof}}(\mathcal{M}^{\text{ret}}(q_i, R_u, k); x_i).$$

Finally, at the prediction stage (Figure 2c), an LLM-based predictor  $\mathcal{M}^{\text{pred}}$  takes post  $x_i$ , comment  $c_i$ , and the user profile  $P_i$  as input to predict whether  $c_i$  will change the user’s view expressed in  $x_i$ :

$$\tilde{y}_i = \mathcal{M}^{\text{pred}}(x_i, c_i; P_i).$$

#### 3.3 Training

Our framework consists of two learnable components: the query generation module ( $\phi^{\text{query}}$ ) and

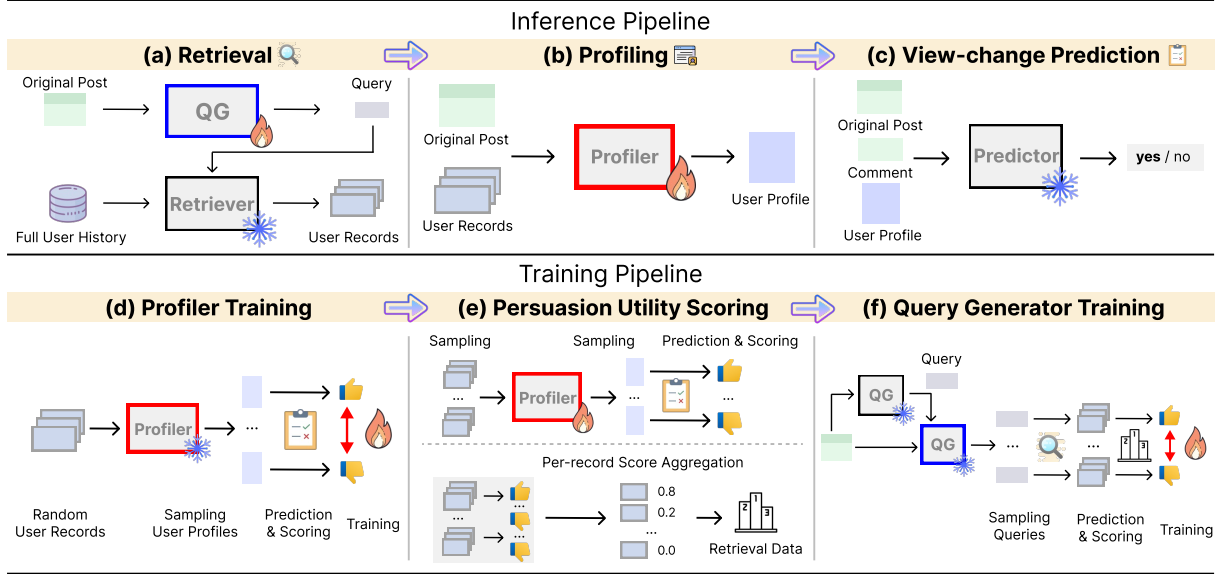


Figure 2: Overview of the proposed framework. *Flame* and *snowflake* icons denote trainable and frozen models, respectively. QG denotes our query generator model. The top illustrates the inference pipeline with three stages: (a) retrieval, (b) profiling, and (c) view-change prediction. The bottom shows the training pipeline, consisting of (d) profiler training, (e) record-level persuasion utility scoring, and (f) query generator training.

the profiler ( $\phi^{\text{prof}}$ ). The overall training pipeline proceeds in three stages: *profiler training*, *persuasion utility scoring*, and *query generator training* (Figure 2d–f). We first train the profiler using randomly retrieved records (Section 3.3.1), and then train the query generator (Section 3.3.3) using supervision derived from persuasion utility scoring (Section 3.3.2), which estimates each record’s contribution to view-change prediction performance.

### 3.3.1 Profiler

Our core hypothesis is that an ideal user profile should be optimized for personalized view-change prediction rather than merely summarizing user history. Since ground-truth profiles for this task are unavailable, we adopt a weakly supervised approach and optimize the profiler via DPO (Figure 2d). We use end-task performance as the preference signal: profiles that successfully predict view change are treated as *chosen*, while unsuccessful ones are treated as *rejected*. To construct such preference data, for each post  $x$ , we randomly sample multiple groups of historical records and prompt the base profiler to generate candidate profiles. We evaluate each profile by its task performance across all comments associated with  $x$ , using the resulting F1 score as a measure of profile quality. Preference pairs are then constructed by pairing higher-scoring profiles with lower-scoring ones separated by a sufficient F1 margin, and the profiler is trained to

prefer higher-quality profiles. Training details are provided in Appendix C.

### 3.3.2 Record-Level Persuasion Utility Scoring

After training the profiler, we estimate the *persuasion utility* of individual user records. This step enables learning in the retrieval stage: training a record-selection module requires supervision that reflects how useful each record is for predicting persuasion outcomes. However, existing datasets lack ground-truth annotations identifying which individual records are most informative for persuasion prediction, and collecting such labels from Reddit users is infeasible. We therefore derive record-level supervision by estimating each record’s contribution to view-change prediction performance (Figure 2e). For each record  $r \in R_u$ , we estimate its contribution by evaluating its effect across different record sets that include  $r$ . Specifically, we randomly partition the records into groups of five, and repeat this grouping process three times. For each group, we generate three user profiles using the trained profiler with a decoding temperature of 0.7. As a result, each record  $r$  is associated with a total of nine generated profiles. We aggregate the F1 scores from all view-change prediction instances performed using profiles that include record  $r$ , and use the result as its persuasion utility score.

### 3.3.3 Query Generator

For effective user profiling, we must retrieve historical user records informative for persuading each user. A naive approach would use the post text directly as a retrieval query. However, CMV posts often lack explicit user-specific attributes critical for persuasion—such as underlying values, relevant experiences, or decision-making styles. Post-only queries thus often fail to retrieve sufficiently useful records for predicting view change. To address this, we train an LLM-based query generator to produce user-focused queries that explicitly target attributes absent from the post but potentially critical for persuasion.

Because training the model to directly infer these implicit attributes from the post alone is difficult, we adopt a two-stage training for the query generator (Figure 2f). First, we prompt the model to generate a user-focused question that targets information not present in the post but likely to influence persuasion (e.g., for a healthcare policy post: “What are the user’s core values regarding government intervention in individual choice?”). Second, we train the model to take both the post and the generated user-focused question as input and generate a single retrieval query that contextualizes the user attribute using salient cues from the post (e.g., “Does the user prioritize individual autonomy over collective benefit when it comes to healthcare access?”). By learning to ground user attributes in the context of the post, the model can better identify which user information is likely to affect persuasion outcomes.

For each post, we sample multiple candidate queries, retrieve user records for each, and score each query by NDCG@5 based on the persuasion utility of the retrieved records. The model is then trained via DPO to prefer queries that yield higher-quality retrieval. At inference, the trained query generator receives only the post and outputs a user-focused query that effectively surfaces persuasion-relevant user records. Full details of candidate generation, preference construction, and optimization are provided in Appendix D.

## 4 Experiments

In this section, we evaluate our proposed framework through two complementary analyses: (1) a retrieval-side evaluation of our query generation strategy based on persuasion utility scores (Section 4.1); (2) end-to-end view-change prediction

Query Strategy	Mean NCG@5	Mean NDCG@5
Random	0.6173	0.6080
BGE-Post	0.6267	0.6180
BGE-Post-Tuned	0.6280	0.6162
HyDE	0.6229	0.6126
<i>Ours</i>	<b>0.6357</b>	<b>0.6214</b>

Table 1: Retrieval performance of different query formulation strategies. Random reports the average performance over 10 runs. BGE-Post and BGE-Post-Tuned use an embedding-based retriever based on BGE-M3, with and without retriever fine-tuning, respectively.

performance, which evaluates the combined effect of all pipeline components (Section 4.2).

**Experimental Setup** We use data collected from the CMV Reddit forum, as described in Section 3.1. Detailed dataset statistics are provided in Appendix A. For the learnable components, we employ Llama-3.1-8B-Instruct as the backbone for both the query generator and the profiler. For embedding-based retrieval, we use the BGE-M3 embedding model. As predictor models, we use two open-weight models at different scales, Llama-3.1-8B-Instruct and Llama-3.3-70B-Instruct, and a closed-source model, GPT-4o-mini, to assess whether our trainable user profiling framework generalizes across predictors. We evaluate performance using the F1 score, which is well-suited to the inherently imbalanced comment labels (e.g., few delta comments among many non-delta comments).

### 4.1 Retrieval-side Experiments

**Setup** We analyze the retrieval component in isolation, focusing on how query generation affects the retrieval of persuasion-relevant user records. Specifically, we compare a random baseline (mean over 10 runs), embedding-based retrieval methods, and our query generation strategies using pre-computed utility scores for individual records (Section 3.3.2). For embedding-based retrieval baselines, we evaluate BGE-POST, which directly use the original post text as the retrieval query, and HYDE (Gao et al., 2023), which generates a hypothetical document that approximates the retrieval target and uses it as the query. Concretely, for HYDE, we prompt Llama-3.1-8B-Instruct with the original post to generate a plausible user record that is likely to be relevant in the given persuasion context.

Method	Llama-3.1-8B-Instruct		Llama-3.3-70B-Instruct		GPT-4o mini	
	F1	AUC	F1	AUC	F1	AUC
No Personalization	0.3457	0.5677	0.3284	0.6538	0.2525	<b>0.6415</b>
PAG	0.2571	0.5775	0.3141	0.6346	0.0833	0.6165
Recursumm	0.3133	0.5869	0.4139	0.6571	0.1050	0.6318
Hsumm	0.3244	0.5965	0.4063	0.6615	0.1128	0.6214
Retrieval-only	0.2952	0.5424	0.4177	0.6635	0.1323	0.6306
<i>Ours</i>	<b>0.4000</b>	<b>0.6158</b>	<b>0.4661</b>	<b>0.6828</b>	<b>0.2787</b>	0.6299

Table 2: End-to-end comparison of our proposed framework with prior user profiling approaches. The table reports F1 and area under the ROC curve (AUC) for view-change prediction across three predictor models.

**Results** Table 1 reports retrieval performance measured by utility-based NCG@5 and NDCG@5. The results indicate that dense retrieval using the post text as the query (BGE-Post) is inherently limited, and that fine-tuning the retriever under the same post-only query formulation (BGE-Post-Tuned) leads to only marginal improvements. This suggests that the bottleneck is not retriever capacity but the incompleteness of the post as a query for eliciting persuasion-relevant user attributes. In contrast, our method improves over BGE-Post and HyDE by transforming the post into a user-focused query that explicitly targets missing attributes conditioned on the post. Consistent with this, end-to-end results (Section 4.2) show that our framework achieves the best view-change prediction performance, highlighting that persuasion-aware query formulation is more beneficial for the full pipeline.

## 4.2 End-to-End View-Change Prediction

**Setup** We evaluate the end-to-end view-change prediction performance of our overall pipeline, comparing it against (1) existing personalized profiling frameworks, and (2) different combinations of retrieval and profiling baselines.

We first compare our method with prior **user profiling frameworks** for personalized dialogue and retrieval-augmented generation, including PAG (Richardson et al., 2023), HSUMM (Zhong et al., 2024), and RECURSUMM (Wang et al., 2025). Details of these baselines are provided in Appendix F. We additionally evaluate two ablations: NO PERSONALIZATION, which performs view-change prediction without user profiles or historical records, and RETRIEVAL-ONLY, which conditions the predictor on raw retrieved records without profile construction.

Next, we conduct a more detailed comparison across different retriever-profiler combinations. For

**retrieval** variants, we compare embedding-based strategies evaluated in Section 4.1 (BGE-POST and HYDE), a sparse retrieval baseline using the post as the query (BM25-POST), and heuristic baselines (RANDOM and RECENT). For **profiling** variants, we consider three approaches to user profile construction: (i) DEMOGRAPHIC, which extracts demographic attributes from retrieved records using GPT-4.1-mini (Hackenburg et al., 2025b; Salvi et al., 2024b), (ii) BASE PROFILER, an instruction-tuned LLM without additional training prompted to summarize retrieved records, and (iii) DPO PROFILER, our profiler trained via DPO (Section 3.3.1).

**Results** Table 2 shows that existing personalization frameworks transfer poorly to view-change prediction. These methods primarily aim to generate user-aligned responses or compress a user’s history, and generic profiles can even hurt performance for Llama-3.1-8B-Instruct and GPT-4o-mini compared to NO PERSONALIZATION. In contrast, our framework yields consistent gains across predictors, achieving a +13.77%p absolute improvement on Llama-3.3-70B-Instruct over NO PERSONALIZATION, indicating that task-oriented, trainable profiling is crucial for personalized persuasion prediction. Compared to the retrieval-only baseline, our approach yields consistent gains across all predictors, ranging from +4.84%p to +14.64%p, highlighting the critical role of the profiler.

Table 3 further decomposes performance by retriever-profiler combinations. Our DPO-trained profiler consistently outperforms demographic and base profiling baselines across all predictors, while demographic profiles perform poorly, suggesting that persuasion-relevant signals are not well captured by demographics alone. On the retrieval side, our query generator delivers the strongest end-to-end performance overall; notably, RECENT is

Retrieval	Llama-3.1-8B-Instruct			Llama-3.3-70B-Instruct			GPT-4o mini		
	Demograph.	Base	Ours	Demograph.	Base	Ours	Demograph.	Base	Ours
Recent	0.3364	0.3805	0.3951	0.3891	<b>0.4058</b>	0.4428	0.0714	0.1629	0.2533
Random	0.3199	0.3758	0.3860	0.4038	0.3979	0.4304	0.0578	0.1516	0.2476
BM25	0.3286	0.3636	0.3742	0.3905	0.3981	0.4218	0.0720	0.1658	0.2754
BGE	0.3286	0.3410	0.3554	<b>0.3912</b>	0.3799	0.4454	0.0663	0.1465	0.2441
HyDE	0.3344	0.3701	0.3785	0.3800	0.3917	0.4507	0.0720	<b>0.1805</b>	0.2570
<i>Ours</i>	<b>0.3466</b>	<b>0.3893</b>	<b>0.4000</b>	0.3837	0.3929	<b>0.4661</b>	<b>0.0765</b>	0.1695	<b>0.2787</b>

Table 3: Effect of retriever and profiler choices on view-change prediction under different predictors (F1). Random reports the average performance over 10 random runs. Underlined results denote our final proposed method, while **boldface** highlights the best-performing configuration within each column. Column groups correspond to different predictor models, with sub-columns indicating profiler configurations (demographic, base profiler, and our trained profiler). Corresponding results using the AUC metric are reported in Appendix G.2.

a competitive baseline, which aligns with Zhang et al. (2025). Our query generator shows substantial synergy with the trained profiler, highlighting that record-level scoring using the trained profiler provides a clear learning signal. Together, these results highlight that view-change prediction benefits most from profiles optimized for the task and retrieval queries that expose persuasion-relevant user attributes, rather than from generic personalization pipelines or standard post-only retrieval.

## 5 Analysis

### 5.1 Profiler Analysis

In this section, we analyze the impact of profiler training by comparing profiles generated by the base profiler (*original profiles*) and the trained profiler (*trained profiles*) on the test set. We present two key analyses below, with results for all predictor models reported in Appendix H.

**(1) The effectiveness of profiler training varies by post topic.** To analyze how profiler effectiveness varies across post characteristics, we annotate each post by *topic* and *claim type* using GPT-4.1-mini. Topics are categorized into *Political* (27.4%), *Sociomoral* (46.4%), and *Others* (26.2%), and claim types into *Interpretation* (39.9%) and *Evaluation* (60.1%), following prior work (Hidey et al., 2017; Priniski and Horne, 2018). Across most topic-claim combinations, trained profiles consistently outperform original profiles in F1. The only exception is political posts under Llama-3.1-8B-Instruct, where profiling benefits sociomoral and other topics but not political posts (Figure 3), likely due to the dominance of group identities in political persuasion. Overall, the results suggest that the trained profiler effectively

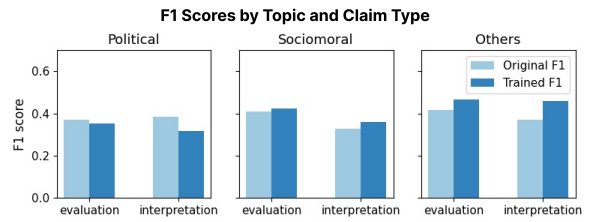


Figure 3: F1 by topic and claim type of the post, comparing the original and trained profilers. Llama-3.1-8B-Instruct is used as the predictor model.

captures individual-level characteristics relevant to persuasion.

**Important profile dimensions vary by post characteristics and the predictor model.** To further analyze profile content, we decompose each profile into sentence-level units, referred to as *profile items*. We annotate each item for the presence of five profile dimensions using GPT-4.1-mini; the five dimensions—*Values & Ideologies*, *Emotional Characteristics*, *Cognitive Characteristics*, *Personality Traits*, and *Interests & Knowledge*—are constructed from persuasion literature (Fabrigar and Petty, 1999; Al Khatib et al., 2020). For each profile, we count the frequency of profile items associated with each dimension and compute *Profile-F1*, the F1 score aggregated over all comments associated with the profile. Figure 4a shows the correlation between (i) change in item frequency for each dimension from original to trained profiles and (ii) the change in Profile-F1.

Our analysis yields three key findings: (1) No single profile dimension is consistently beneficial or detrimental across all posts. (2) The effect of each dimension is strongly post-dependent: for example, cognitive traits (e.g., reasoning styles, de-

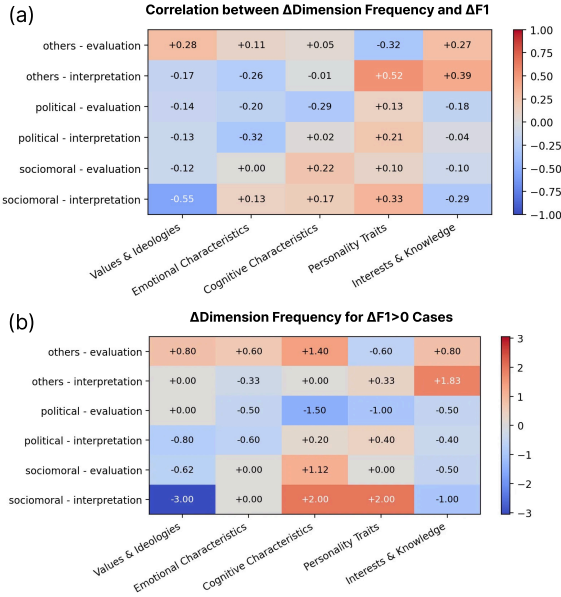


Figure 4: Analysis of profile-dimension frequency shifts ( $\Delta DF$ ) and performance gains ( $\Delta F1$ ) between the original and trained profilers. (a) Correlation between  $\Delta DF$  and  $\Delta F1$ . (b)  $\Delta DF$  for cases with  $\Delta F1 > 0$ . Llama-3.1-8B-Instruct is used as the predictor model.

cision making styles) are positively associated with performance gains for political evaluation posts, but negatively associated with sociomoral evaluation posts. (3) For cases where performance improves, shifts in item frequency across profile dimensions align with these correlation trends (Figure 4b), indicating that different dimensions are emphasized depending on the post. These three patterns remain consistent across different predictor models. However, the specific association patterns between post characteristics and profile dimensions vary substantially with the choice of predictor model. Taken together, these results suggest that persuasion-relevant dimensions differ across posts and the predictor models, and that the trained profiler captures this post-dependent and predictor-specific variation by dynamically adjusting the emphasis on different dimensions. Results for all predictor models are reported in Appendix H.

## 5.2 User Record Analysis

We conduct an analysis of user records scored via persuasion-utility scoring, using LLAMA-3.1-8B-INSTRUCT as the predictor model, focusing on (1) semantic differences between high- and low-scoring records and (2) cross-model patterns in utility scores. More detailed analyses are provided

Model Pair	Top-5 Overlap	Spearman $\rho$
GPT / Llama70B	0.273	0.007
GPT / Llama8B	0.281	0.083
Llama70B / Llama8B	0.245	-0.005

Table 4: Pairwise agreement between predictors on record-level utility scores, measured by mean top-5 overlap and Spearman  $\rho$ . GPT, Llama8B, and Llama70B correspond to GPT-4o-mini, Llama-3.1-8B-Instruct, and Llama-3.3-70B-Instruct, respectively.

in Appendix I.

**Low-scoring records are not semantically dissimilar to the post.** We analyze pairs of top-1 and bottom-1 records for the same post in the validation set, focusing on cases where the bottom-1 record receives an F1 score of zero. Following Section 5.1, we annotate the records along *topic* and *claim type*. Contrary to the hypothesis that low-scoring records fail due to semantic misalignment with the post, we observe the opposite trend: low-scoring records are more likely than high-scoring ones to share the same topic or claim type as the post. This highlights the need for finer-grained contextualization in persuasion.

**Different predictor models prefer different user records.** We further compare persuasion utility scores across predictor models and find little agreement in their preferred records (Table 4). Pairwise comparisons of the top-5 records show low overlap (0.24–0.28), corresponding to only about 1.25 shared records on average. Similarly, Spearman rank correlations are near zero (-0.005–0.083), indicating weak consistency in relative ordering. These results indicate that record utility for view-change prediction is highly model-dependent, motivating the training of predictor-specific retrieval modules.

## 6 Conclusion

We introduce a trainable user profiling framework that captures persuasion-relevant user factors. Experiments on the CMV dataset show that our approach consistently outperforms baselines by constructing context-dependent profiles tailored to the downstream predictor. By learning to retrieve and construct task-oriented user profiles, our framework enables scalable, context-sensitive personalization without retraining predictors or requiring extensive user annotations, making it practical for real-world decision-support systems such as conversational agents, recommendation, and coaching.

## 605 **Limitations**

606 This study focuses on personalized persuasive-  
607 ness prediction in the setting of online opin-  
608 ion change, evaluated on the ChangeMyView Red-  
609 dit dataset. While this setting provides a well-  
610 established testbed with explicit view-change sig-  
611 nals, it represents a specific form of persuasion  
612 grounded in long-form textual discussions. Extend-  
613 ing the framework to other interaction modalities  
614 or domains—such as short-form conversations or  
615 real-time recommendation settings—would require  
616 additional validation.

## 617 **Ethical Considerations**

618 Research on predicting view change in online dis-  
619 cussions could be related to ethical considerations  
620 about user autonomy and the responsible use of  
621 predictive insights. In this work, we strictly focus  
622 on predicting whether a view change occurs, rather  
623 than intervening in user behavior. The framework  
624 does not include any mechanisms for generating  
625 persuasive content or intervening on individuals,  
626 but is designed to enhance understanding of view-  
627 change dynamics in natural settings.

628 All experiments are conducted using publicly  
629 available and anonymized data, without any per-  
630 sonally identifiable information.

## 631 **AI Assistance Acknowledgement**

632 We used AI assistants to proofread the writing and  
633 to help with coding.

## References

- 635 Khalid Al Khatib, Michael Völske, Shahbaz Syed, Niko-  
636 lay Kolyada, and Benno Stein. 2020. Exploiting  
637 personal characteristics of debaters for predicting  
638 persuasiveness. In *Proceedings of the 58th Annual  
639 Meeting of the Association for Computational Lin-  
640 guistics*, pages 7067–7072.
- 641 Esin Durmus and Claire Cardie. 2018a. Exploring the  
642 role of prior beliefs for argument persuasion. In *Pro-  
643 ceedings of the 2018 Conference of the North Amer-  
644 ican Chapter of the Association for Computational  
645 Linguistics: Human Language Technologies, Volume  
646 1 (Long Papers)*, pages 1035–1045, New Orleans,  
647 LA. Association for Computational Linguistics.
- 648 Esin Durmus and Claire Cardie. 2018b. [Exploring the  
649 role of prior beliefs for argument persuasion](#). In *Pro-  
650 ceedings of the 2018 Conference of the North Amer-  
651 ican Chapter of the Association for Computational  
652 Linguistics: Human Language Technologies, Volume  
653 1 (Long Papers)*, pages 1035–1045, New Orleans,  
654 Louisiana. Association for Computational Linguistics.  
655
- 656 Esin Durmus and Claire Cardie. 2019a. [A corpus for  
657 modeling user and language effects in argumentation  
658 on online debating](#). In *Proceedings of the 57th An-  
659 nual Meeting of the Association for Computational  
660 Linguistics*, pages 602–607, Florence, Italy. Associa-  
661 tion for Computational Linguistics.
- 662 Esin Durmus and Claire Cardie. 2019b. [Modeling the  
663 factors of user success in online debate](#). In *The World  
664 Wide Web Conference (WWW) 2019*, pages 2701–  
665 2707.
- 666 Leandre R Fabrigar and Richard E Petty. 1999. The  
667 role of the affective and cognitive bases of attitudes  
668 in susceptibility to affectively and cognitively based  
669 persuasion. *Personality and social psychology bul-  
670 letin*, 25(3):363–381.
- 671 Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan.  
672 2023. [Precise zero-shot dense retrieval without rel-  
673 evance labels](#). In *Proceedings of the 61st Annual  
674 Meeting of the Association for Computational Lin-  
675 guistics (Volume 1: Long Papers)*, pages 1762–1777,  
676 Toronto, Canada. Association for Computational Lin-  
677 guistics.
- 678 Amirata Ghorbani and James Zou. 2019. [Data shapley:  
679 Equitable valuation of data for machine learning](#). In  
680 *Proceedings of the 36th International Conference on  
681 Machine Learning (ICML)*, pages 2242–2251.
- 682 Kobi Hackenburg, Ben M. Tappin, Luke Hewitt,  
683 Ed Saunders, Sid Black, Hause Lin, Catherine Fist,  
684 Helen Margetts, David G. Rand, and Christopher  
685 Summerfield. 2025a. The levers of political persua-  
686 sion with conversational artificial intelligence. *Sci-  
687 ence*, 390(6777):eaea3884.
- 688 Kobi Hackenburg, Ben M Tappin, Luke Hewitt,  
689 Ed Saunders, Sid Black, Hause Lin, Catherine Fist,  
Helen Margetts, David G Rand, and Christopher Sum-  
merfield. 2025b. The levers of political persuasion  
with conversational artificial intelligence. *Science*,  
390(6777):eaea3884.
- Christopher Hidey and Kathleen McKeown. 2018. [Per-  
suaasive influence detection: The role of argument  
sequencing](#). In *Proceedings of the AAAI Conference  
on Artificial Intelligence (AAAI)*.
- Christopher Hidey, Elena Musi, Alyssa Hwang,  
Smaranda Muresan, and Kathleen McKeown. 2017.  
Analyzing the semantic types of claims and premises  
in an online persuasive forum. In *Proceedings of the  
4th Workshop on Argument Mining*, pages 11–21.
- Lu Ji, Zhongyu Wei, Xiangkun Hu, Yang Liu, Qi Zhang,  
and Xuanjing Huang. 2018. [Incorporating argument-  
level interactions for persuasion comments evaluation  
using co-attention model](#). In *Proceedings of the 27th  
International Conference on Computational Linguis-  
tics*, pages 3703–3714, Santa Fe, New Mexico, USA.  
Association for Computational Linguistics.
- Khalid Al Khatib, Michael Völske, Shahbaz Syed, Niko-  
lay Kolyada, and Benno Stein. 2020. Exploiting  
personal characteristics of debaters for predicting  
persuasiveness. In *Proceedings of the 58th Annual  
Meeting of the Association for Computational Lin-  
guistics*, pages 7067–7072. Association for Compu-  
tational Linguistics.
- Cheng Li, Mingyang Zhang, Qiaozhu Mei, Weize Kong,  
and Michael Bendersky. 2024. Learning to rewrite  
prompts for personalized text generation. In *Pro-  
ceedings of the ACM Web Conference 2024*, pages  
3367–3378.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Sp-  
ithourakis, Jianfeng Gao, and Bill Dolan. 2016. [A  
persona-based neural conversation model](#). In *Pro-  
ceedings of the 54th Annual Meeting of the Associa-  
tion for Computational Linguistics (Volume 1: Long  
Papers)*, pages 994–1003, Berlin, Germany. Associa-  
tion for Computational Linguistics.
- Stephanie Lukin, Pranav Anand, Marilyn Walker, and  
Steve Whittaker. 2017a. Argument strength is in the  
eye of the beholder: Audience effects in persuasion.  
In *Proceedings of the 15th Conference of the Euro-  
pean Chapter of the Association for Computational  
Linguistics (EACL)*, pages 742–753.
- Stephanie Lukin, Pranav Anand, Marilyn Walker, and  
Steve Whittaker. 2017b. Argument strength is in the  
eye of the beholder: Audience effects in persuasion.  
In *Proceedings of the 15th Conference of the Euro-  
pean Chapter of the Association for Computational  
Linguistics: Volume 1, Long Papers*, pages 742–753.
- Sheshera Mysore, Zhuoran Lu, Mengting Wan, Longqi  
Yang, Bahareh Sarrafzadeh, Steve Menezes, Tina  
Baghaee, Emmanuel Barajas Gonzalez, Jennifer  
Neville, and Tara Safavi. 2024. Pearl: Personal-  
izing large language model writing assistants with  
generation-calibrated retrievers. In *Proceedings of*



854	into training, validation, and test sets with an ap-	You will then be shown a post you wrote,	900
855	proximate ratio of 8:1:1. Table 5 presents the de-	and a comment written in response to it.	901
856	tailed statistics for each split, including the distribu-	Based on your history, determine whether	902
857	tion of user history length and the volume of delta	the comment would change your mind from	903
858	and non-delta comments.	the opinion expressed in the post.	904
859	<b>B Predictor Model Prompts</b>	Respond only with one word: "yes" if	905
860	In this section, we provide the detailed prompts	your mind would change after reading the	906
861	used for the predictor models. We present the Sys-	comment, or "no" if not. Do not provide	907
862	tem Prompt and User Prompt sequentially for each	any explanation or reasoning.	908
863	setting.	<b>User Prompt</b>	909
864	<b>B.1 Prediction with User Profile (Ours)</b>	### User History	910
865	<b>System Prompt</b>	{user_profile}	911
866	You are the author of the post. The	### Post	912
867	section labeled "User Profile" is your	{post}	913
868	profile – it describes who you are.	### Comment	914
869	Read it carefully and fully adopt this	{comment}	915
870	as your identity and mindset.	–	916
871	You will then be shown a post you wrote,	Would this comment change your mind from	917
872	and a comment written in response to it.	the opinion you expressed in the post?	918
873	Based on your profile, determine whether	Respond only with one word: "yes" or	919
874	the comment would change your mind from	"no".	920
875	the opinion expressed in the post.	<b>B.3 Prediction without Personalization</b>	921
876	Respond only with one word: "yes" if	<b>System Prompt</b>	922
877	your mind would change after reading the	You are the author of the post. Carefully	923
878	comment, or "no" if not. Do not provide	read your own post and the comment	924
879	any explanation or reasoning.	written in response to it.	925
880	<b>User Prompt</b>	Decide whether you would change your	926
881	### User Profile	mind after reading the comment.	927
882	{user_profile}	Ignore your own beliefs as a language	928
883	### Post	model and fully adopt the mindset of the	929
884	{post}	person who wrote the post.	930
885	### Comment	Respond with only one word: "yes" if you	931
886	{comment}	think you would change your mind, or "no"	932
887	–	if not. Do not provide any explanation	933
888	Would this comment change your mind from	or reasoning.	934
889	the opinion you expressed in the post?	<b>User Prompt</b>	935
890	Respond only with one word: "yes" or	[Post]	936
891	"no".	{post}	937
892	<b>B.2 Prediction with User History</b>	[Comment]	938
893	<b>(Retrieval-Only)</b>	{comment}	939
894	<b>System Prompt</b>	Would you change your mind after reading	940
895	You are the author of the post.	the comment?	941
896	The section labeled "User History" is	<b>C Profiler Training Details</b>	942
897	relevant past history about you.	In this section, we provide detailed specifications	943
898	Read it carefully and incorporate it	for the preference construction process and the hy-	944
899	into your identity and mindset.	perparameters used for Direct Preference Optimiza-	945
		tion (DPO).	946

Split	# Posts	Unique OPs	User History Count				Delta Comments			Non-Delta Comments		
			Min	Max	Mean	Median	Min	Mean	Median	Min	Mean	Median
Train	1,341	1,257	15	11,965	252.40	57	1	1.77	1	1	33.06	20.0
Validation	167	69	16	19,583	956.35	65	1	2.24	1	1	31.81	19.0
Test	168	69	15	19,583	613.36	71	1	2.54	1.5	2	35.30	19.0

Table 5: Detailed statistics of the dataset splits. User History Count refers to the number of historical posts/comments available for the OP prior to the current post.

### C.1 Preference Pair Construction

To derive robust training signals from the synthesized candidate profiles, we employ a margin-based stratified sampling strategy. As described in Section 3.3.1, for each input group  $\mathcal{G}_i$ , we generate a set of 16 candidate profiles  $\Pi_{\mathcal{G}_i}$ . We rank these profiles based on their utility score  $S(\pi)$ , which represents the macro-F1 score on the view-change prediction task.

To avoid noisy training signals arising from pairs with negligible performance differences, we enforce a minimum utility margin  $\delta$ . We construct a dataset of preference pairs  $\mathcal{D} = \{(x, \pi_w, \pi_l)\}$  where:

$$S(\pi_w) - S(\pi_l) \geq \delta \quad (1)$$

where  $x$  represents the input historical records. Specifically, we select the top- $K$  performing profiles as positive samples and the bottom- $K$  profiles as negative samples from the candidate set  $\Pi_{\mathcal{G}_i}$ . We then form pairs from the Cartesian product of these two subsets, filtering out any pairs that do not satisfy the margin condition. In our experiments, we set  $K = 4$  (top 25% and bottom 25%) and the margin  $\delta = 0.05$  to ensure distinct quality separation.

### C.2 DPO Training Configuration

We optimize the profiler  $\pi_\theta$  using the standard DPO objective, which increases the likelihood of the preferred profile  $\pi_w$  while decreasing that of the dispreferred profile  $\pi_l$ , implicitly optimizing the reward function without a separate reward model training step. The loss function is defined as:

$$\begin{aligned} \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = & \\ -E_{(x, \pi_w, \pi_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(\pi_w|x)}{\pi_{\text{ref}}(\pi_w|x)} \right) \right. & (2) \\ \left. - \beta \log \frac{\pi_\theta(\pi_l|x)}{\pi_{\text{ref}}(\pi_l|x)} \right) & \end{aligned}$$

where  $\pi_{\text{ref}}$  is the frozen reference model (the initial base profiler),  $\sigma$  is the logistic sigmoid function,

and  $\beta$  is a hyperparameter controlling the deviation from the reference model.

We initialized the profiler with Llama-3.1-8B-Instruct. To ensure training stability and prevent overfitting to the small number of high-utility patterns, we utilized Low-Rank Adaptation (LoRA) for parameter-efficient fine-tuning. The detailed hyperparameters are listed in Table 6.

Hyperparameter	Value
Base Model	Llama-3.1-8B-Instruct
LoRA Rank ( $r$ )	32
LoRA Alpha ( $\alpha$ )	64
Optimizer	AdamW
Learning Rate	5e-7
LR Scheduler	Linear
Warmup Ratio	0.05
Batch Size	64
Beta ( $\beta$ )	0.1
Epochs	3
Max Sequence Length	16384

Table 6: DPO training hyperparameters for the profiler.

### C.3 Profile Generation Prompts

We use the following prompts to generate a context-aware user profile tailored for persuasion.

#### System Prompt

You are an expert assistant whose task is to extract concise, high-level information about the author of a set of passages.

Focus only on traits that would be most useful for persuading or changing the author’s view in relation to the current post.

Your goal is to produce a compact, context-aware user profile optimized for persuasive messaging toward the given post.

#### User Prompt

You are given a set of passages written by the same author, along with the author’s current post.

1011	Extract only the most essential	persuasion-critical user attributes (e.g., values, ex-	1062
1012	information about the author that	periences, decision-making styles). To make learn-	1063
1013	is clearly stated or strongly and	ing easier, we adopt a two-stage training strategy:	1064
1014	consistently implied across multiple		
1015	passages, focusing on traits that are	1. <b>Stage 1 (User-Focused Question Genera-</b>	1065
1016	most relevant for understanding how to	<b>tion).</b> Prompt the model to produce a user-	1066
1017	persuade them in the context of the	focused <i>question</i> that asks for user informa-	1067
1018	current post.	tion <i>not present in the post</i> but likely to affect	1068
1019	Instructions:	persuasion.	1069
1020	- Consider the current post as the	2. <b>Stage 2 (Post-Contextualized Query Gener-</b>	1070
1021	immediate context in which persuasion	<b>ation).</b> Train the model to take the post and	1071
1022	would occur.	the Stage-1 question as input and generate a	1072
1023	- Identify attitudes, reasoning	single <i>retrieval query</i> that contextualizes the	1073
1024	patterns, or sensitivities that could	user attribute using salient post cues (topic,	1074
1025	influence how the author might respond	stance, constraints).	1075
1026	to persuasion regarding the post.		
1027	- Do not guess or speculate beyond what	In second stage, supervision is derived from	1076
1028	is well supported.	retrieval quality: we score candidate queries by	1077
1029	- Exclude personally identifying	NDCG@5 based on the persuasion utility of re-	1078
1030	or sensitive details unless explicitly	trieved records and apply DPO to prefer candi-	1079
1031	stated.	dates with higher retrieval quality. At inference,	1080
1032	- Generalize from specific events or	the trained model receives only the post and out-	1081
1033	examples into higher-level traits; avoid	puts a user-focused retrieval query.	1082
1034	direct quotes or low-level details.		
1035	- Remove redundancy and keep bullets	<b>D.2 Candidate Generation Procedure</b>	1083
1036	concise.	For each post $x_i$ , we generate candidates as fol-	1084
1037	- Do NOT respond with anything other	lows.	1085
1038	than the bullet points.		
1039	Current Post:	<b>Stage 1: User-Focused Question.</b> We first gen-	1086
1040	{post}	erate a single user-focused question $q_i^{(1)}$ from the	1087
1041	Input Passages:	query generator using the Stage-1 prompt with de-	1088
1042	{passages}	coding temperature 0. This question serves as an	1089
1043	Output:	intermediate representation of the user attribute to	1090
1044	• ...	seek.	1091
1045	• ...		
1046		<b>Stage 2: Post-Contextualized Retrieval Query.</b>	1092
1047	This appendix provides implementation details for	Conditioned on $(x_i, q_i^{(1)})$ , we sample 16 candi-	1093
1048	training the query generator, including candidate	date retrieval queries $\{q_{i,j}^{(2)}\}_{j=1}^{16}$ using the Stage-2	1094
1049	generation, retrieval-based supervision, preference	prompt with temperature 0.8. Each candidate is a	1095
1050	construction, and optimization settings. Across	single natural-language sentence that integrates (i)	1096
1051	all experiments, the query generator is imple-	the user attribute targeted by $q_i^{(1)}$ and (ii) salient	1097
1052	mented as a single LLM (Llama-3.1-8B-Instruct)	cues from $x_i$ .	1098
1053	and trained using Direct Preference Optimization	This two-step candidate generation is used to	1099
1054	(DPO) with a two-stage training strategy, following	construct DPO training data and is applied consis-	1100
1055	Section 3.3.3.	tently across all predictor models.	1101
1056			
1057	<b>D.1 Overview</b>	<b>D.3 Retrieval and Scoring</b>	1102
1058	The query generator is trained to produce <i>user-</i>	Each Stage-2 candidate query $q_{i,j}^{(2)}$ is used to re-	1103
1059	<i>focused retrieval queries</i> that retrieve historical	trieve the top-5 user records from the author’s his-	1104
1060	user records informative for personalized persua-	torical records using a fixed embedding-based re-	1105
1061	sion. A post-only query is often insufficient	triever (BGE-M3). We evaluate query quality using	1106
	because the post may not explicitly mention		

Predictor	Pos.	Neg.	Max
Llama-3.1-8B-Instruct	$\geq 0.65$	$\leq 0.55$	8
Llama-3.3-70B-Instruct	$\geq 0.75$	$\leq 0.65$	8
GPT-4o-mini	$\geq 0.55$	$\leq 0.45$	10

Table 7: Predictor-specific thresholds for preference pair construction based on NDCG@5.

Hyperparameter	Value
Base model	Llama-3.1-8B-Instruct
LoRA rank $r$	16
LoRA scaling $\alpha$	32
Learning rate	$2 \times 10^{-5}$
Max epochs	3
DPO $\beta$	0.3 (0.1 for GPT-4o-mini)

Table 8: DPO training hyperparameters for the query generator.

NDCG@5, where the graded relevance of each retrieved record is given by its pre-computed *record-level persuasion utility score* (Section 3.3.2).

Both the retriever and the utility scores are kept fixed throughout query generator training. Thus, the query generator is trained solely to improve retrieval quality under a fixed downstream evaluation signal.

#### D.4 Preference Pair Construction

For each post, we partition the 16 Stage-2 candidates into positive and negative pools based on their NDCG@5 scores, and construct preference pairs by pairing positives with negatives. We additionally enforce a minimum margin of 0.10 between the chosen and rejected query scores, and select up to a fixed maximum number of pairs per post.

Because the NDCG@5 score distributions differ across predictor models (due to predictor-specific persuasion utility scoring), we use predictor-specific thresholds and pair caps to ensure sufficient supervision. Table 7 summarizes the settings.

#### D.5 Optimization via DPO

We train the query generator using DPO with LoRA fine-tuning. All settings are shared across predictors except the DPO inverse temperature  $\beta$ , which we tune to account for differences in preference sharpness under each predictor’s utility signal.

#### D.6 Query Generator Prompts

##### User-Focused Question Prompt.

##### System Prompt

You will be given an online post where a user explains their view on a specific topic.

Write ONE short question that asks for information regarding the user that is NOT explicitly stated in the post, but would be important for persuading the user expressed in the post.  
The question should focus on aspects such as the user’s values, experiences, priorities, or decision making styles related to the topic.

Instructions:

- Output MUST be a single question sentence ending with “?”.
- Do NOT explain your reasoning.
- Do NOT ask for information already provided in the post.

##### User Prompt

Post:

–  
{post}

Respond in ONE question.

##### Post-Contextualized Query Prompt.

##### System Prompt

You will be given two inputs:

- (1) an online post where a user explains their view on a specific topic.
- (2) a question asking for information that is NOT explicitly stated in the post, but is important for persuading the user in this situation.

Write ONE sentence that incorporates:

- what the question is asking about the user
- the most important cues from the post

The sentence should clearly reflect what the question asks about the user, while also grounding it in the most important cues from the post.

Instructions:

- Output MUST be a single sentence.
- Do NOT explain your reasoning.

##### User Prompt

Post:

–  
{post}

Question:

–  
{question}

Respond in ONE sentence.

- The author values practicality and efficiency in decision-making, prioritizing the enjoyment and benefit gained from an experience over the cost.
- The author is likely a rational and logical thinker, weighing the pros and cons of a situation before making a decision.
- The author may be skeptical of the idea of "getting one's money's worth" and instead focuses on the value gained from an experience.
- The author values personal enjoyment and satisfaction over material possessions or financial gain.
- The author may be open to re-evaluating past decisions and choices if they no longer align with their current values or interests.
- The author is likely resistant to the idea of doing something solely for the sake of "getting their money's worth" if it doesn't bring them enjoyment or benefit.

- The author values clarity and predictability in laws, seeking a system where laws are clearly defined and easily understandable.
- They believe in the importance of accountability, particularly for those in power, and want to reduce the potential for bias and corruption.
- The author is concerned about the potential for abuse of power and the need for checks and balances.
- They are skeptical of the current system's ability to hold those in power accountable, particularly when it comes to lawmakers and the government.
- The author is interested in exploring alternative solutions to address issues such as corruption and conflict of interest.
- They value the idea of a more autonomous system, where the court system is reserved for complex or unclear cases.
- The author is open to considering the possibility of starting a court case for an unrelated crime, especially if there is strong evidence and the authorities are refusing to take action.
- They are likely to be influenced by logical and evidence-based arguments, particularly those that highlight the potential flaws in the current system.

- Values self-restraint and responsibility, particularly in the context of societal well-being.
- Recognizes the importance of considering the potential consequences of individual actions on others.
- Is skeptical of simplistic or absolute advice, such as "be yourself," and prefers nuanced thinking.
- Understands that personal flaws and imperfections can have negative impacts on others.
- Is willing to acknowledge and work on personal limitations, rather than pretending to be perfect.
- May be concerned with the potential negative outcomes of certain actions, such as rape and murder, and sees these as self-damning.
- May be open to considering the potential consequences of individual actions on the greater good.
- May be interested in exploring the complexities of human behavior and the factors that influence it.
- May be critical of absolute or dogmatic thinking.

Figure 5: Examples of user profiles generated by our trained profiler when used Llama3.1 8B Instruct as the predictor model for training.

## E Examples of the Generated User Profiles

To provide qualitative insight into the behavior of our trained profiler, we present examples of generated user profiles at Figure 5.

Figure 6 compares profiles generated for the same user when trained with different predictor models. Despite being grounded in identical user records, the resulting profiles exhibit systematic differences in emphasis and framing, reflecting predictor-specific preferences for information that is most useful for downstream prediction. This comparison highlights that our profiler does learn to adapt profile construction tailored to the target predictor model.

## F Details of User Profiling Baselines

We compare our approach against several representative user profiling frameworks. Specifically, we consider: (i) PAG (Richardson et al., 2023), which

### Llama 3.1 8B Instruct

- Values self-restraint and responsible behavior, particularly in the context of societal impact.
- Believes that individual flaws and impulses can have negative consequences for others.
- Prioritizes the greater good and the well-being of society over personal desires or impulses.
- Is skeptical of simplistic or overly idealistic advice, such as "be yourself," when it may lead to negative outcomes.
- Recognizes the importance of personal growth and self-improvement, as evident in the desire to "become a good citizen."
- May be concerned with the potential consequences of unchecked individual behavior on the social fabric.
- May be open to nuanced and context-dependent advice, rather than one-size-fits-all solutions.

### Llama 3.3 70B Instruct

- Values self-restraint and responsibility, particularly in the context of societal well-being.
- Recognizes the importance of considering the potential consequences of individual actions on others.
- Is skeptical of simplistic or absolute advice, such as "be yourself," and prefers nuanced thinking.
- Understands that personal flaws and imperfections can have negative impacts on others.
- Is willing to acknowledge and work on personal limitations, rather than pretending to be perfect.
- May be concerned with the potential negative outcomes of certain actions, such as rape and murder, and sees these as self-damning.
- May be open to considering the potential consequences of individual actions on the greater good.
- May be interested in exploring the complexities of human behavior and the factors that influence it.
- May be critical of absolute or dogmatic thinking.

### GPT 4o mini

- The author is likely open to critical thinking and nuanced perspectives, as they acknowledge the limitations of a simplistic "be yourself" approach and consider various counterexamples.
- They value pragmatism and the greater good, as they prioritize becoming a "good citizen" over personal flaws and their potential consequences.
- The author is likely skeptical of absolute moral or philosophical principles, as they question the idea that people would behave poorly without a higher power.
- They may be concerned with social responsibility and the impact of individual actions on society, as they consider the potential consequences of unrestrained behavior.
- The author is willing to admit their own flaws and imperfections, indicating a level of self-awareness and humility.
- They may be interested in exploring the complexities of human nature and the role of societal factors in shaping individual behavior.
- The author is likely open to considering multiple perspectives and evaluating evidence, as they engage in a debate and provide counterexamples to a given argument.

Figure 6: Examples of user profiles generated by our trained profilers for the same user under different predictor models used for training.

retrieves a subset of user records using BM25, independently summarizes each selected record, and concatenates the resulting summaries into a user profile; (ii) HSUMM (Zhong et al., 2024), which applies hierarchical summarization by first generating summaries over subsets of user records and then aggregating these intermediate summaries into a single profile; (iii) RECURSUMM (Wang et al., 2025), which incrementally updates the user profile by recursively integrating each new user record with the existing summary. PAG relies on retrieval to select a subset of user records for profiling, whereas HSUMM, RECURSUMM, and PRIME assume access to and utilize all available user historical records when constructing user profiles. Concretely, HSUMM and RECURSUMM construct a profile for each user by repeatedly summarizing that user's entire history through multiple summa-

1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228

Query Strategy	Llama-3.3-70B-Instruct		GPT-4o-mini	
	NCG@5	NDCG@5	NCG@5	NDCG@5
Random	0.7461	0.7395	0.4713	0.4647
BGE-Post	0.7461	0.7357	0.4826	0.4736
HyDE	0.7528	<b>0.7482</b>	0.4685	0.4562
<i>Ours</i>	<b>0.7536</b>	0.7471	<b>0.4827</b>	<b>0.4747</b>

Table 9: Additional retrieval-side results under different predictor-specific persuasion utility signals. All methods are evaluated on the same record pool as the main experiment. Random reports the average over 10 runs.

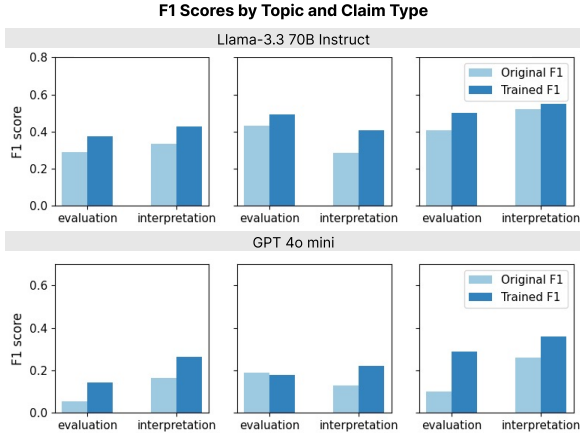


Figure 7: F1 by topic and claim type of the post, comparing the original and trained profilers.

1229 rization steps.

## 1230 G Additional Experiment Results

### 1231 G.1 Additional Retrieval Results

1232 Table 9 reports retrieval performance under dif- 1252  
 1233 ferent predictor-specific persuasion utility signals, 1253  
 1234 using the same candidate record pool and query 1254  
 1235 strategies as in the main experiment (Table 1).

### 1236 G.2 AUC scores across different retrieval and 1255 1237 profiling variants 1256

1238 Table 10 reports the full end-to-end AUC results 1257  
 1239 across different retriever and profiler combina- 1258  
 1240 tions for each predictor model. 1259  
 1241

## 1242 H Details of Profiler Analysis

### 1243 H.1 Additional Results

1244 Figure 7 extends the main-text analysis by re- 1252  
 1245 porting F1 scores by topic and claim type 1253  
 1246 for GPT-4o-mini and Llama-3.3-70B-Instruct, 1254  
 1247 comparing the original and trained profilers. 1255

1248 Figures 8 and 9 present the profile- 1256  
 1249 dimension analysis results for GPT-4o-mini 1257  
 and Llama-3.3-70B-Instruct, respectively. 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264

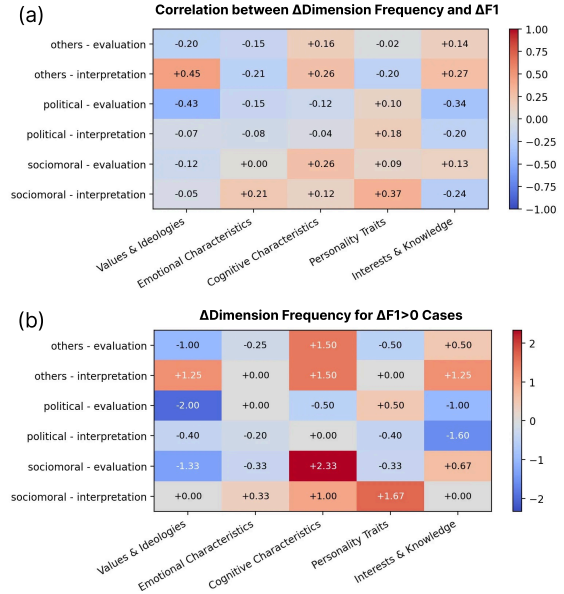


Figure 8: Analysis of profile-dimension frequency shifts ( $\Delta DF$ ) and performance gains ( $\Delta F1$ ) between the original and trained profilers. (a) Correlation between  $\Delta DF$  and  $\Delta F1$ . (b)  $\Delta DF$  for cases with  $\Delta F1 > 0$ . GPT-4o-mini is used as the predictor.

## I Details of User Record Analysis

### I.1 Additional Results

Figure 10 analyzes the topical and claim type characteristics of user records ranked by utility score. Top-ranked records tend to align more closely with the post in both topic and claim type, whereas bottom-ranked records show weaker alignment.

Table 11 reports mean F1 over all records as well as the top-5 and bottom-5 subsets, showing that while Llama-3.3-70B-Instruct assigns higher scores overall, GPT-4o-mini exhibits a larger contrast between high- and low-ranked records. This trend is further reflected in Table 12, where GPT-4o-mini yields consistently larger margins between top-5 and bottom-5 records.

Retrieval	Llama-3.1-8B-Instruct			Llama-3.3-70B-Instruct			GPT-4o-mini		
	Demograph.	Base	Ours	Demograph.	Base	Ours	Demograph.	Base	Ours
Recent	0.5828	0.5953	0.6088	0.6577	0.6740	0.6746	0.6214	0.6214	0.6232
Random	<b>0.5859</b>	<b>0.6121</b>	0.6112	0.6528	0.6669	0.6716	0.6188	0.6121	0.6253
BM25	0.5858	0.5955	0.6082	0.6564	0.6588	0.6697	0.6159	0.6189	0.6365
BGE	0.5851	0.5859	0.6029	<b>0.6596</b>	<b>0.6768</b>	0.6798	0.6226	0.6305	0.6349
HyDE	0.5845	0.5997	0.6104	0.6569	0.6655	0.6825	0.6216	<b>0.6311</b>	<b>0.6447</b>
<i>Ours</i>	0.5850	0.6054	<b>0.6146</b>	0.6574	0.6736	<b>0.6828</b>	<b>0.6232</b>	0.6020	<u>0.6299</u>

Table 10: Effect of retriever and profiler choices on view-change prediction under different predictors (AUC). Random reports the average performance over 10 runs. Underlined results denote our final proposed method, while **boldface** highlights the best-performing configuration within each column. Column groups correspond to different predictor models, with sub-columns indicating profiler configurations (demographic, base profiler, and our trained profiler).

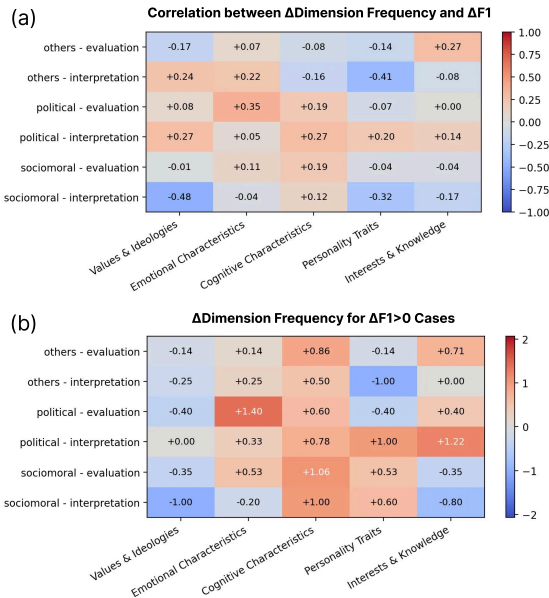


Figure 9: Analysis of profile-dimension frequency shifts ( $\Delta$ DF) and performance gains ( $\Delta$ F1) between the original and trained profilers. (a) Correlation between  $\Delta$ DF and  $\Delta$ F1. (b)  $\Delta$ DF for cases with  $\Delta$ F1 > 0. Llama-3.3-70B-Instruct is used as the predictor.

Model	All	Top-5	Bottom-5
GPT-4o-mini	0.163	0.302	0.092
Llama-3.3-70B-Instruct	<b>0.375</b>	<b>0.468</b>	<b>0.301</b>
Llama-3.1-8B-Instruct	0.248	0.321	0.183

Table 11: Mean F1 scores across all records, top-5 records, and bottom-5 records. While Llama-3.3-70B-Instruct assigns higher absolute scores overall, GPT-4o-mini exhibits a larger contrast between top-5 and bottom-5 records.

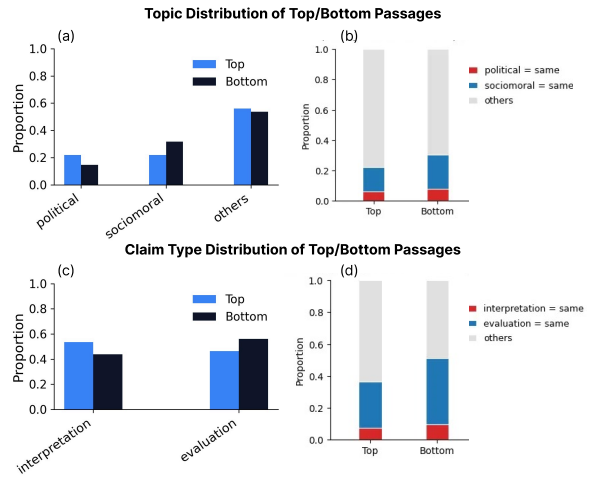


Figure 10: Topic and claim type distributions of top and bottom-ranked records. (a,c) Marginal distributions by topic and claim type. (b,d) Proportions of records sharing the same topic or claim type with the post.

Model	Mean	Median (p50)	p90
GPT-4o-mini	<b>0.180</b>	<b>0.133</b>	<b>0.500</b>
Llama3.3-70B-Instruct	0.135	0.090	0.337
Llama-3.1-8B-Instruct	0.112	0.080	0.260

Table 12: Distribution of margin between high- and low-scoring records (min top-5 minus max bottom-5). GPT-4o-mini exhibits consistently larger margins, indicating clearer separation between beneficial and non-beneficial records.