

# STEP-BACK PROFILING: Distilling User Interactions for Personalized Scientific Writing

Anonymous Authors

## Abstract

Large language models (LLMs) excel at a variety of natural language processing tasks, yet they struggle to generate personalized content for individuals, particularly in real-world settings like scientific writing. Addressing this challenge, we introduce STEP-BACK PROFILING that personalizes LLMs by abstracting user interactions into concise profiles. Our approach effectively condenses user interaction history, distilling it into profiles that encapsulate essential traits and preferences of users, thus facilitating personalization that is both effective and user-specific. Importantly, STEP-BACK PROFILING is a low-cost and easy-to-implement technique that does not require additional fine-tuning. Through evaluation of the LaMP benchmark, which encompasses a spectrum of language tasks requiring personalization, our approach outperformed the baseline, showing improvements of up to 3.6 points. We curated the Personalized Scientific Writing (PSW) dataset to further study multi-user personalization in challenging real-world scenarios. This dataset requires the models to write scientific papers given specialized author groups with diverse academic backgrounds. On PSW, we demonstrate the value of capturing collective user characteristics via STEP-BACK PROFILING for collaborative writing. Extensive experiments and analysis validate our method’s state-of-the-art performance and broader applicability – an advance that paves the way for more user-tailored scientific applications with LLMs.

## 1 Introduction

In recent years, Large Language Models (LLMs) have made significant strides in natural language understanding and generation, demonstrating human-parity performance on a wide range of tasks [Wei *et al.*, 2022b,a; Chowdhery *et al.*, 2023; OpenAI, 2023]. Moreover, the advent of LLM-driven language agents has revolutionized a myriad of user-facing applications, marking a game-changing breakthrough in the general AI capacity [Zhou *et al.*, 2023; Zhang *et al.*, 2023b; Shinn *et al.*, 2023; Qin *et al.*, 2023; Qiao *et al.*, 2024]. Con-

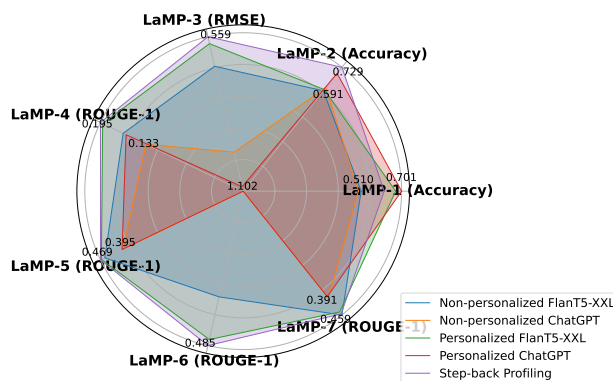


Figure 1: STEP-BACK PROFILING consistently improves the downstream task accuracy on the LaMP dataset.

currently, integrating LLMs with personalization paradigms has paved the way for a vast frontier in improving user-centric services and applications [Salemi *et al.*, 2023; Chen *et al.*, 2023; Zhiyuli *et al.*, 2023], as they provide a deeper understanding of users’ accurate demands and interests than abstract vector-based information representations. By learning to characterize and emulate user-specific language patterns, personalized LLMs can enable more engaging and valuable interactions in domains such as dialogue [Wang *et al.*, 2019; Zhang *et al.*, 2019b; Character.AI, 2022], recommendation [Zhiyuli *et al.*, 2023; Wang *et al.*, 2023], role-playing [Shao *et al.*, 2023; Jiang *et al.*, 2023] and content creation [Cao *et al.*, 2023; Wei *et al.*, 2022c].

Prior work on personalizing language models [Salemi *et al.*, 2023; Tan and Jiang, 2023; Zhang *et al.*, 2023a; Chen *et al.*, 2023; Zhiyuli *et al.*, 2023] has shown promise, but primarily focused on learning user representations in a single-user context. For example, the LaMP benchmark [Salemi *et al.*, 2023] evaluates personalization given a single target user’s historical interactions on tasks like citation prediction and product review generation. However, many real-world applications involve multiple users collaborating on a shared task, such as team-authored scientific papers.

Another practical challenge for LLM personalization is scaling to extensive user histories while respecting context

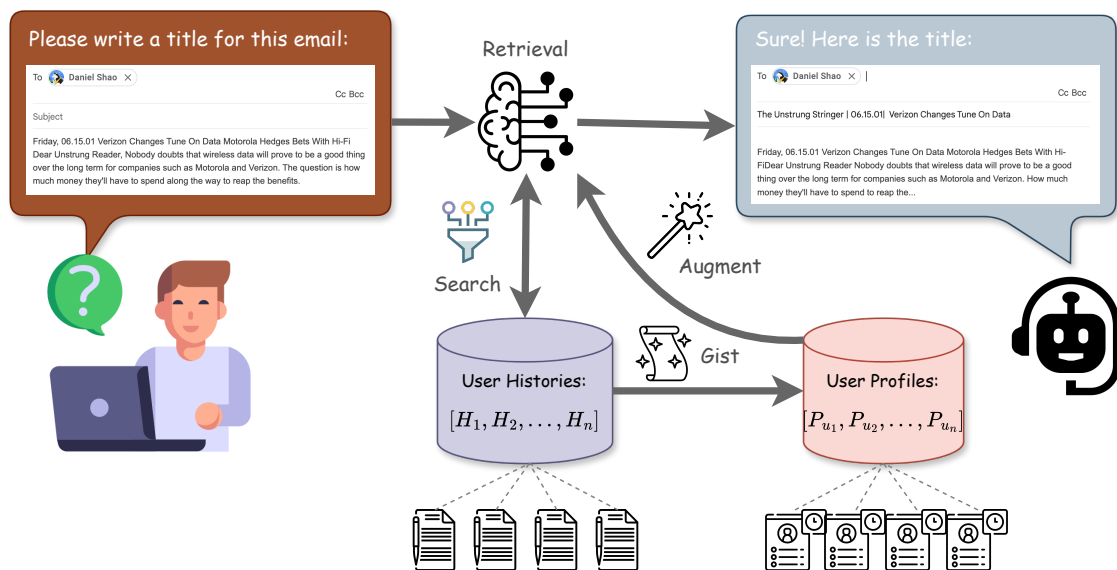


Figure 2: Overview of the STEP-BACK PROFILING Methodology. It applies the ‘gist’ abstraction function to the history of users and generates personalized output through an enhanced retrieval-augmented language model.

67 length limits [Shi *et al.*, 2023; Liu *et al.*, 2024]. Directly  
 68 conditioning on raw personal histories quickly becomes in-  
 69 feasible as user data grows. Prior methods mostly use un-  
 70 compressed history for personalization [Salemi *et al.*, 2023],  
 71 which restricts the amount of user-specific information the  
 72 model can utilize. This limits knowledge-intensive applica-  
 73 tions like scientific writing, where relevant information may  
 74 be diluted across many documents.

75 This work proposes a training-free LLM personalization  
 76 framework that addresses these challenges through STEP-  
 77 BACK PROFILING – inspired by the ideas of gist memory  
 78 [Lee *et al.*, 2024] and STEP-BACK PROMPTING [Zheng *et al.*,  
 79 2023] for information compression and abstraction, we distill  
 80 individual user histories into concise profile representations  
 81 that capture high-level concepts and language traits. This en-  
 82 ables efficient memory management and allows the model  
 83 to focus on salient user characteristics, grounding personal-  
 84 ized generation without excess computation or laborious data  
 85 collection [Chen *et al.*, 2023]. STEP-BACK PROFILING is a  
 86 low-cost and easy-to-implement technique that operates di-  
 87 rectly on the pre-trained LLM without additional training. It  
 88 can also complement parameter-efficient techniques of LLM  
 89 finetuning [Hu *et al.*, 2021; Dettmers *et al.*, 2024; Sheng *et al.*,  
 90 2023]. We show that STEP-BACK PROFILING improves per-  
 91 formance over standard personalization methods on the  
 92 LaMP benchmark.

93 Moreover, we introduce the Personalized Scientific Writ-  
 94 ing (PSW) dataset to study multi-user personalization. PSW  
 95 contains research papers collaboratively written by expert  
 96 teams, and each author’s background publications are used  
 97 to construct profiles. Modeling a group’s collective expertise  
 98 is crucial for this task, as different paper sections may re-  
 99 flect knowledge associated with particular authors. PSW thus  
 100 poses a challenging and realistic testbed for multi-user per-  
 101 sonalization, requiring both abstractions of individual exper-

tise and dynamic integration of diverse user traits throughout  
 the collaborative writing process.

To summarize the contributions of this work:

1. A training-free STEP-BACK PROFILING approach that  
 enables efficient and expressive personalization by ab-  
 stracting user histories into trait-centric representations.
2. The Personalized Scientific Writing (PSW) dataset, a  
 real-world benchmark for studying multi-user personal-  
 ization with a novel task of collaborative expert writing.
3. A state-of-the-art performance of STEP-BACK PROFIL-  
 ING for single and multi-user personalization on diverse  
 tasks in the LaMP and PSW benchmark without addi-  
 tional training.

## 2 STEP-BACK PROFILING

### 2.1 Motivation

Existing methods for personalizing language models struggle  
 to effectively utilize user histories, particularly in the pres-  
 ence of extraneous details that can obscure the most pertinent  
 information for a given task [Shi *et al.*, 2023; Liu *et al.*, 2024].  
 This challenge is magnified in multi-user scenarios, where  
 models must efficiently extract and integrate knowledge from  
 multiple users’ histories. While retrieval-augmented meth-  
 ods, such as those employed in the LaMP benchmark [Salemi  
*et al.*, 2023], have made progress in scaling to more extensive  
 user histories, they still operate on raw user data containing  
 relevant and irrelevant details. To address these limitations,  
 we introduce a STEP-BACK PROFILING approach that dis-  
 tills a user’s raw history into a concise representation focusing  
 on ‘gist’ representations and preferences, drawing inspira-  
 tion from the STEP-BACK PROMPTING technique [Zheng *et al.*,  
 2023] and the READAGENT framework [Lee *et al.*, 2024].  
 Our approach aims to enable more efficient and effective per-  
 sonalization across diverse single and multi-user scenarios by

135 reasoning about higher-level traits instead of verbatim user  
136 history.

## 137 2.2 Procedure

138 Consider a set of  $n$  users denoted by  $U = \{u_1, u_2, \dots, u_n\}$ ,  
139 where each user  $u_i$  has a preference history  $H_i =$   
140  $\{(x_{i1}, y_{i1}), (x_{i2}, y_{i2}), \dots, (x_{im}, y_{im})\}$  consisting of  $m$   
141 input-output pairs. To effectively generate  $P(y|x, H_U)$  based  
142 on users' preference history, we create a set of user profiles  
143  $P_U = \{P_{u_1}, P_{u_2}, \dots, P_{u_n}\}$  using STEP-BACK PROFILING.  
144 The complete procedure involves the following steps:

145 **User Profile Gisting:** Each user's history is condensed into  
146 a short "gist" representation using an abstraction function  
147  $\text{Gist}(\cdot)$ :  $P_{u_i} = \text{Gist}(H_i)$ . The "gist" captures the user's  
148 high-level traits and interests.

149 **Multi-User Profile Concatenation:** Individual user pro-  
150 files  $\{P_{u_1}, P_{u_2}, \dots, P_{u_n}\}$  are concatenated to form a unified  
151 representation  $P_U$ :  $P_U = [P_{u_1}; P_{u_2}; \dots; P_{u_n}]$ , where  $[\cdot; \cdot]$  is  
152 a permutation-sensitive function combining the user profiles.

153 **Retrieval-Augmented Generation (Optional):** Relevant  
154 snippets from user histories  $H_U$  may be retrieved for  
155 input  $x$  using a retrieval function  $\text{Retrieve}(\cdot)$ :  $R_i =$   
156  $\text{Retrieve}(x, H_i, k)$ , where  $R_i$  is a set of top- $k$  retrieved input-  
157 output snippets from user  $u_i$ 's history  $H_i$ . The retrieved snip-  
158 pets  $R = \{R_1, R_2, \dots, R_n\}$  can be concatenated with  $x$  to  
159 form an augmented input  $\hat{x}$ :  $\hat{x} = [x; R_1; R_2; \dots; R_n]$ .

160 **Personalized Output Generation:** The personalized lan-  
161 guage model generates an output  $y$  by conditioning on the  
162 augmented input  $\hat{x}$  (if retrieval is used) or the original in-  
163 put  $x$ , along with the concatenated user profile  $P_U$ :  $y =$   
164  $\text{Generate}(\hat{x}, P_U)$ . The generated output  $y$  aligns with the  
165 user preferences captured by the STEP-BACK PROFILING  
166 while following the input  $x$ .

## 167 3 The Personalized Scientific Writing (PSW) 168 Benchmark

169 We have extended the LaMP benchmark, introduced by  
170 Salemi *et al.* [2023], to evaluate multi-user scenarios. Our  
171 PSW benchmark includes four tasks, and we outline the data  
172 collection process for each task.

### 173 3.1 Problem Formulation

174 Personalized language models aim to generate outputs that  
175 follow a given input and align with the users' styles, prefer-  
176 ences, and expertise. In multi-author collaborative writing,  
177 the Personalized Writing Styles (PSW) benchmark provides  
178 a framework for evaluating such models.

179 Each data entry in the PSW benchmark consists of four key  
180 components:

- 181 1. An input sequence  $x$  serves as the model's input.
- 182 2. A target output  $y$  that the model is expected to generate.
- 183 3. A set of user histories  $H_U = \{H_{u_1}, H_{u_2}, \dots, H_{u_k}\}$ ,  
184 where  $k$  is the number of collaborating authors, and each  
185 entry  $H_{u_i}$  contains historical input-output pairs for user  
186  $u_i$ .

4. A set of author roles  $C = \{c_1, c_2, \dots, c_k\}$ , each rep- 187  
resenting the role of the corresponding author  $u_i$  in the 188  
collaborative writing process. 189

190 A personalized language model aims to generate an out-  
191 put  $y$  that aligns with the conditional probability distribution  
192  $P(y|x, H_U, C)$ . This means the model should produce an  
193 output that follows the input  $x$  and the collaborating authors'  
194 writing styles, preferences, and expertise, as captured by their  
195 user histories  $H_U$  and roles  $C$ .

196 By conditioning the language model's output on these ad-  
197 ditional factors, the PSW benchmark allows for the evaluation  
198 of personalized models that can adapt to the unique charac-  
199 teristics of individual authors in a collaborative writing envi-  
200 ronment.

### 201 3.2 Task Descriptions

202 **UP-0: Research Interest Generation:** Before all the PSW  
203 tasks, we create a benchmark for user profiling. This in-  
204 volves compiling a list of research interests that accurately  
205 reflect each author's expertise and research focus based on  
206 their publication history. To acquire the necessary informa-  
207 tion, we extract the research interests of each author from  
208 Google Scholar<sup>1</sup> by searching their name. Once we have this  
209 information, we incorporate it into the author's profile.

210 **PSW-1: Research Topic Generation:** This task aims to  
211 generate a list of research topics that capture the collaborat-  
212 ing authors' joint expertise and research focus, given their  
213 user profiles. The generated research topics should be rele-  
214 vant to the authors' past publications and help identify poten-  
215 tial research directions for their collaborative work. We use  
216 OpenAI's *gpt-4* model to automatically extract research top-  
217 ics from selected papers. The extracted topics are then linked  
218 to their respective papers and author profiles.

219 **PSW-2: Research Question Generation:** This task fo-  
220 cuses on generating a set of research questions that align with  
221 the expertise and interests of the collaborating authors and are  
222 relevant to the target paper. The generated research questions  
223 should help guide the content and structure of the collabora-  
224 tive writing process. We automatically use OpenAI's *gpt-4*  
225 model to extract research questions from the selected papers  
226 for this task. The extracted research questions are then linked  
227 to their papers and author profiles.

228 **PSW-3: Paper Abstract Generation:** This task involves  
229 generating a paper abstract that summarizes the key points  
230 and contributions of the collaborative research paper, given  
231 the user profiles, research interests, target paper title, and re-  
232 search questions. The generated abstract should incorporate  
233 the writing styles and preferences of the collaborating authors  
234 while maintaining coherence and clarity. For this task, we  
235 directly retrieve the abstracts from the selected papers using  
236 the Semantic Scholar API<sup>2</sup>. The retrieved abstracts are then  
237 linked to their respective papers and author profiles.

238 **PSW-4: Paper Title Generation:** This task aims to gener-  
239 ate a suitable title for the collaborative research paper, consid-  
240 ering the user profiles, research interests, research questions,

<sup>1</sup><https://github.com/scholarly-python-package/scholarly>

<sup>2</sup><https://api.semanticscholar.org/>

241 and paper abstract. The generated title should be concise,  
 242 informative, and reflect the paper’s main contributions while  
 243 considering the collaborating authors’ expertise and interests.  
 244 The data for this task is collected using the Semantic Scholar  
 245 API, which provides the titles of the selected papers.

### 246 3.3 G-Eval for PSW Evaluation

247 We use the G-Eval framework [Liu *et al.*, 2023] to evaluate  
 248 the generated outputs on the PSW benchmark. G-Eval em-  
 249 ploys LLMs like GPT-4 with chain-of-thought prompting to  
 250 assess the quality of generated text in a form-filling paradigm  
 251 [Zhang *et al.*, 2019a]. G-Eval is particularly well-suited for  
 252 evaluating the PSW tasks because it can handle open-ended  
 253 generation tasks without requiring gold reference outputs and  
 254 provides scores that closely approximate expert human judg-  
 255 ments [Yuan *et al.*, 2021]. We can use the G-Eval frame-  
 256 work to obtain multi-dimensional evaluations of PSW model  
 257 outputs. These dimensions include consistency, fluency, rel-  
 258 evance, and novelty, which are considered essential scientific  
 259 writing criteria [Kryściński *et al.*, 2019; Fabbri *et al.*, 2021].  
 260 An example G-Eval prompt can be found in Appendix C.

## 261 4 Experimental Setup

262 We assess our methods alongside other baseline approaches  
 263 across the LaMP and PSW datasets. This section provides  
 264 a detailed exploration of the experimental settings for these  
 265 evaluations.

### 266 4.1 Datasets and Evaluation

267 **LaMP Dataset:** We follow the standard practice estab-  
 268 lished in Salemi *et al.* [2023], encompassing three classifi-  
 269 cation and four text generation tasks. Specifically, these tasks  
 270 are Personalized Citation Identification (LaMP-1), Personal-  
 271 ized News Categorization (LaMP-2), Personalized Product  
 272 Rating (LaMP-3), Personalized News Headline Generation  
 273 (LaMP-4), Personalized Scholarly Title Generation (LaMP-  
 274 5), Personalized Email Subject Generation (LaMP-6), and  
 275 Personalized Tweet Paraphrasing (LaMP-7).

276 **PSW Dataset:** As introduced in the previous section, the  
 277 PSW dataset is designed to assess the performance of per-  
 278 sonalized language models in collaborative scientific writing  
 279 scenarios. The dataset includes one individual task, User Pro-  
 280 filing (UP-0), and four collaborative tasks: Research Topics  
 281 Generation (PSW-1), Research Question Generation (PSW-  
 282 2), Paper Abstract Generation (PSW-3), and Paper Title Gen-  
 283 eration (PSW-4).

284 **Evaluation:** Our evaluation methodology mirrors the  
 285 LaMP framework outlined in Salemi *et al.* [2023]. We  
 286 evaluate our proposed methods using the metrics specified  
 287 in the LaMP benchmark for each task. These include F1  
 288 score, Accuracy, MAE, and RMSE for classification tasks and  
 289 ROUGE-1 and ROUGE-L for generation tasks.

### 290 4.2 Methods to compare

291 We employ *gpt-3.5-turbo* hosted by OpenAI for all tasks in  
 292 this paper. Our proposed method is compared against sev-  
 293 eral baselines, including non-personalized language models,

models fine-tuned on history data without personalization,  
 and models that use simple concatenation of user histories  
 for personalization.

## 297 4.3 Main Result

Dataset	Metric	Non-personalized		Personalized		STEP-BACK PROFILING
		FlanT5-XXL	ChatGPT	FlanT5-XXL	ChatGPT	
LaMP-1	Accuracy	0.522	0.510	0.675	<b>0.701</b>	0.624
LaMP-2	Accuracy	0.591	0.610	0.598	0.693	<b>0.729</b>
	F1	0.463	0.455	0.477	0.455	<b>0.591</b>
LaMP-3	MAE	0.357	0.699	0.282	0.658	<b>0.274</b>
	RMSE	0.666	0.977	0.584	1.102	<b>0.559</b>
LaMP-4	ROUGE-1	0.164	0.133	0.192	0.160	<b>0.195</b>
	ROUGE-L	0.149	0.118	0.178	0.142	<b>0.180</b>
LaMP-5	ROUGE-1	0.455	0.395	0.467	0.398	<b>0.469</b>
	ROUGE-L	0.410	0.334	0.424	0.336	<b>0.426</b>
LaMP-6	ROUGE-1	0.332	-	0.466	-	<b>0.485</b>
	ROUGE-L	0.320	-	0.453	-	<b>0.464</b>
LaMP-7	ROUGE-1	<b>0.459</b>	0.396	0.448	0.391	0.455
	ROUGE-L	<b>0.404</b>	0.337	0.396	0.324	0.398

Table 1: Performance comparison of personalized and non-personalized models on the LaMP dataset.

**LaMP Results:** To guarantee a fair comparison, we uti-  
 lized a user-based separation from Salemi *et al.* [2023]. We  
 only granted the agent access to the provided user history  
 and restricted it from accessing any other information. Ad-  
 ditionally, we utilized the same pre-trained retriever, without  
 any additional fine-tuning, to retrieve the top five examples.  
 This approach was identical to the *Non-Personalized* setting  
 in Salemi *et al.* [2023]. Finally, we compared our results with  
 the outcomes reported in the study.

As shown in Table 1<sup>3</sup>, our analysis unveiled a notable  
 performance enhancement through our method’s application,  
 significantly when leveraging the same backbone language  
 models (*gpt-3.5-turbo*). In the domain of text generation tasks  
 (LaMP-4~7), our method achieved an average improvement  
 of 0.048 in Rouge-1 and 0.053 in Rouge-L, corresponding  
 to gains of 15.2% and 19.5%, respectively. Similarly, for  
 the classification tasks (LaMP-1~3), we observed an average  
 +12.6% accuracy gain of and a +42.5% reduction in MAE  
 compared to the *Non-Personalized* setting. Our method con-  
 tinues to exhibit better performance across most tasks, even  
 when compared with **FlanT5-XXL**, with a fine-tuned re-  
 triever as *Personalized* setting. The prompt used in this ex-  
 periment is detailed in Appendix D.

### PSW Results

We assess the proficiency of our proposed personalized  
 agent using the PSW dataset, focusing on user profiling (UP-  
 0), personalized idea brainstorming (PSW-1, PSW-2), and  
 personalized text generation (PSW-3, PSW-4). We compare  
 the performance of our method in three different settings:

1. **Zero-shot:** Generates outputs based on the input  
 prompt  $x$  alone:

$$y = \text{Generate}(x).$$

<sup>3</sup>Baseline results are obtained directly from Salemi *et al.* [2023].

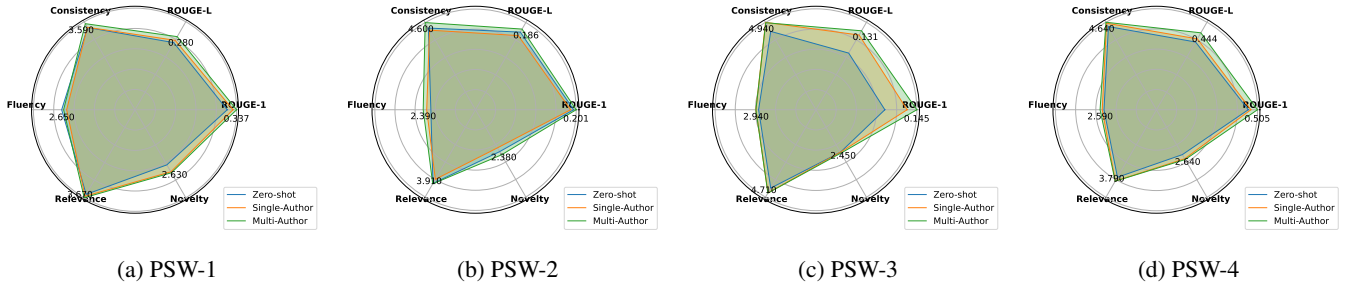


Figure 3: PSW datasets’ Performance metrics (ROUGE-1, ROUGE-L, Consistency, Fluency, Relevance, and Novelty) across three different models: **Zero-shot**, **Single-Author**, and **Multi-Author**. The **Multi-Author** model consistently achieves the highest scores across all datasets.

Datasets Method	Metrics						
	ROUGE-1	ROUGE-L	Consistency	Fluency	Relevance	Novelty	
UP-0 <b>Single-Author</b>	0.267	0.233	4.32	2.01	3.59	/	
PSW-1 <b>Zero-shot</b>	0.306	0.257	3.43	<b>2.65</b>	3.53	2.30	
PSW-1 <b>Single-Author</b>	0.325	0.266	3.44	2.47	3.61	2.59	
PSW-1 <b>Multi-Author</b>	<b>0.337</b>	<b>0.280</b>	<b>3.59</b>	2.58	<b>3.67</b>	<b>2.63</b>	
PSW-2 <b>Zero-shot</b>	0.196	0.179	4.31	2.04	3.89	2.21	
PSW-2 <b>Single-Author</b>	0.190	0.171	4.20	2.23	3.67	2.01	
PSW-2 <b>Multi-Author</b>	<b>0.201</b>	<b>0.186</b>	<b>4.60</b>	<b>2.39</b>	<b>3.91</b>	<b>2.38</b>	
PSW-3 <b>Zero-shot</b>	0.099	0.094	4.43	2.81	4.43	2.40	
PSW-3 <b>Single-Author</b>	0.131	0.124	<b>4.94</b>	<b>2.94</b>	4.70	2.40	
PSW-3 <b>Multi-Author</b>	<b>0.145</b>	<b>0.131</b>	4.92	<b>2.94</b>	<b>4.71</b>	<b>2.45</b>	
PSW-4 <b>Zero-shot</b>	0.459	0.391	4.41	2.41	3.58	2.38	
PSW-4 <b>Single-Author</b>	0.472	0.409	4.59	2.49	3.78	2.60	
PSW-4 <b>Multi-Author</b>	<b>0.505</b>	<b>0.444</b>	<b>4.64</b>	<b>2.59</b>	<b>3.79</b>	<b>2.64</b>	

Table 2: Performance comparison of personalized models on the PSW dataset, with additional metrics such as **Consistency (1-5)**, **Fluency (1-3)**, **Relevance (1-5)**, and **Novelty (1-3)** reported.

+5.1%, +6.7%, +3.8%, and +6.4%, respectively, compared to the **Zero-shot** and **Single-Author** setting. The prompt used in this experiment is detailed in Appendix E.

#### 4.4 Ablation Studies

To assess the contribution of each component, we perform an ablation study on the PSW dataset. Table 3 and 4 report the results of two variants: 1) Switching the order of users and 2) Removing user profiling.

##### Impact of Author Order

Table 3 shows how changing the author order affects the performance of multi-user personalized models. We experiment with three variants:

- **Original**: The original author order as provided in the dataset.
- **Swap-Random**: Randomly shuffle the order of authors.
- **Swap-First**: Move the first author to the end of the author list.

Datasets Variants	Metrics						
	ROUGE-1	ROUGE-L	Consistency	Fluency	Relevance	Novelty	
PSW-1 <b>Original</b>	<b>0.337</b>	<b>0.280</b>	<b>3.59</b>	<b>2.58</b>	3.67	<b>2.63</b>	
PSW-1 <b>Swap-Random</b>	0.321	0.272	3.42	2.48	<b>3.69</b>	2.45	
PSW-1 <b>Swap-First</b>	0.314	0.260	3.35	2.42	3.48	2.37	
PSW-2 <b>Original</b>	<b>0.201</b>	<b>0.186</b>	<b>4.60</b>	<b>2.39</b>	<b>3.91</b>	2.38	
PSW-2 <b>Swap-Random</b>	0.193	0.178	4.53	2.30	3.85	<b>2.42</b>	
PSW-2 <b>Swap-First</b>	0.186	0.171	4.46	2.27	3.77	2.29	
PSW-3 <b>Original</b>	<b>0.145</b>	<b>0.131</b>	<b>4.92</b>	2.94	<b>4.71</b>	<b>2.45</b>	
PSW-3 <b>Swap-Random</b>	0.138	0.125	4.84	2.88	4.65	2.50	
PSW-3 <b>Swap-First</b>	0.130	0.117	4.78	<b>2.98</b>	4.57	2.55	
PSW-4 <b>Original</b>	<b>0.505</b>	<b>0.444</b>	<b>4.64</b>	<b>2.59</b>	<b>3.79</b>	2.64	
PSW-4 <b>Swap-Random</b>	0.492	0.431	4.57	2.55	3.72	2.70	
PSW-4 <b>Swap-First</b>	0.483	0.421	4.50	2.50	3.64	<b>2.76</b>	

Table 3: Impact of author order on the performance of multi-user personalized models, with additional metrics such as **Consistency (1-5)**, **Fluency (1-3)**, **Relevance (1-5)**, and **Novelty (1-3)** reported.

The **Original** order consistently achieves the best performance across all metrics on all PSW tasks. Randomly swapping authors (**Swap-Random**) leads to a slight decline, while moving the first author to the end (**Swap-First**) results in a more significant drop. This observation highlights the importance of preserving the original author order

2. **Single-Author**: Personalizes with single user’s profile  $P_{u_i}$  and retrieved snippets  $R_i$ :

$$y = \text{Generate}(\hat{x}, P_{u_i}),$$

where  $\hat{x} = [x; R_i]$  and  $R_i = \text{Retrieve}(x, H_i, 10)$ .

3. **Multi-Author**: Personalizes with multiple users’ profiles  $P_U$  and retrieved snippets  $R$ :

$$y = \text{Generate}(\hat{x}, P_U),$$

where  $\hat{x} = [x; R_1; \dots; R_n]$ ,  $R_i = \text{Retrieve}(x, H_i, 10)$  for each user  $u_i$ .

As shown in Table 2, our **Multi-Author** setting demonstrates superior performance across all tasks. In the personalized idea brainstorming tasks (PSW-1 and PSW-2), the **Multi-Author** setting outperforms both **Zero-shot** and **Single-Author** settings, with an average improvement of +6.9% in ROUGE-1 and +7.1% in ROUGE-L. Similarly, for the personalized text generation tasks (PSW-3 and PSW-4), the **Multi-Author** setting achieves the highest ROUGE scores, with an average gain of +28.2% in ROUGE-1 and +26.6% in ROUGE-L, compared to the **Zero-shot** and **Single-Author** settings. Furthermore, the **Multi-Author** setting exhibits the highest scores for additional metrics such as Consistency, Fluency, Relevance, and Novelty across all tasks, with an average improvement of

373 in multi-author collaborative writing scenarios. The first au-  
 374 thor, often the lead or corresponding author, significantly in-  
 375 fluences the document’s content, structure, and style. As a  
 376 result, their writing style and expertise tend to be most promi-  
 377 nently reflected in the document. Disrupting this order intro-  
 378 duces noise and hinders the model’s ability to capture the in-  
 379 dividual authors’ impact and the logical progression of ideas,  
 380 particularly affecting the generation tasks (PSW-3 and PSW-  
 381 4) where content and style are heavily influenced by the main  
 382 author’s expertise and preferences.

### 383 Impact of User Profiling

384 Table 4 reports ablation results on the user profile component:

- 385 • **Original:** User profiles constructed using STEP-  
 386 BACK PROFILING.
- 387 • **Removed:** No user profiles used, only retrieving rele-  
 388 vant snippets.
- 389 • **Random:** Replacing target user profiles with randomly  
 390 sampled user profiles.

Datasets Profile		Metrics					
		ROUGE-1	ROUGE-L	Consistency	Fluency	Relevance	Novelty
PSW-1	<b>Original</b>	<b>0.337</b>	<b>0.280</b>	<b>3.59</b>	<b>2.58</b>	<b>3.67</b>	<b>2.63</b>
	<b>Removed</b>	0.297	0.250	3.21	2.49	3.31	2.57
	<b>Random</b>	0.328	0.272	3.55	2.56	3.62	2.68
PSW-2	<b>Original</b>	<b>0.201</b>	<b>0.186</b>	<b>4.60</b>	2.39	<b>3.91</b>	2.38
	<b>Removed</b>	0.180	0.166	4.28	2.32	3.63	2.33
	<b>Random</b>	0.195	0.182	4.57	<b>2.42</b>	3.89	<b>2.45</b>
PSW-3	<b>Original</b>	<b>0.145</b>	<b>0.131</b>	4.92	2.94	<b>4.71</b>	2.45
	<b>Removed</b>	0.128	0.115	4.70	2.87	4.50	2.41
	<b>Random</b>	0.142	0.128	<b>4.95</b>	<b>2.96</b>	4.69	<b>2.51</b>
PSW-4	<b>Original</b>	<b>0.505</b>	<b>0.444</b>	<b>4.64</b>	<b>2.59</b>	<b>3.79</b>	2.64
	<b>Removed</b>	0.475	0.419	4.38	2.53	3.58	2.56
	<b>Random</b>	0.498	0.438	4.60	2.58	3.76	<b>2.69</b>

Table 4: Impact of user profile on the performance of multi-user personalized models, with additional metrics such as **Consistency (1-5)**, **Fluency (1-3)**, **Relevance (1-5)**, and **Novelty (1-3)** reported.

391 Removing user profiles (**Removed**) leads to the largest  
 392 performance decline, confirming the benefit of STEP-BACK  
 393 PROFILING in multi-user personalization. Using random pro-  
 394 file texts (**Random**) recovers some of the gaps but still under-  
 395 performs the **Original** profiles. This demonstrates that the  
 396 distilled user traits successfully capture useful information  
 397 for collaborative writing, such as individual writing styles,  
 398 expertise, and preferences. The performance gap between  
 399 **Original** and **Random** profiles highlights the effective-  
 400 ness of the STEP-BACK PROFILING technique in extracting  
 401 relevant user characteristics from their background informa-  
 402 tion. These findings underscore the importance of incorpo-  
 403 rating author-specific traits to enable a more personalized and  
 404 contextually appropriate generation in multi-user settings.

## 405 5 Conclusion

406 In summary, STEP-BACK PROFILING offers a promising way  
 407 to improve the effectiveness and scalability of personalized  
 408 language models. Abstracting user histories into compact  
 409 profiles enables the model to better focus on pertinent in-  
 410 formation and handle longer contexts. Experiments on both

single-user and multi-user settings validate the benefits of  
 profile-guided personalization.

Future work can explore more advanced profiling strate-  
 gies, such as hierarchical representations and dynamic profile  
 updates based on user feedback. Adapting STEP-BACK PRO-  
 FILING to long histories spanning multiple sessions is another  
 valuable direction. Finally, studying the interpretability and  
 controllability of profile-guided models can help build user  
 trust and allow for more fine-grained customization.

411  
412  
413  
414  
415  
416  
417  
418  
419

## 420 References

- 421 Yihan Cao, Siyu Li, Yixin Liu, Zhiling Yan, Yutong Dai,  
422 Philip S Yu, and Lichao Sun. A comprehensive survey of  
423 AI-generated content (AIGC): A history of generative AI  
424 from GAN to ChatGPT. *arXiv preprint arXiv:2303.04226*,  
425 2023.
- 426 Character.AI. Character.AI. <https://character.ai/>, 2022.
- 427 Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu,  
428 Gangwei Jiang, Yuanhao Pu, Yuxuan Lei, Xiaolong Chen,  
429 Xingmei Wang, et al. When large language models meet  
430 personalization: Perspectives of challenges and opportuni-  
431 ties. *arXiv preprint arXiv:2307.16376*, 2023.
- 432 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,  
433 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul  
434 Barham, Hyung Won Chung, Charles Sutton, Sebastian  
435 Gehrmann, et al. PaLM: Scaling language modeling  
436 with pathways. *Journal of Machine Learning Research*,  
437 24(240):1–113, 2023.
- 438 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke  
439 Zettlemoyer. QLoRA: Efficient finetuning of quantized  
440 LLMs. *Advances in Neural Information Processing Sys-*  
441 *tems*, 36, 2024.
- 442 Alexander R Fabbri, Wojciech Kryściński, Bryan Mc-  
443 Cann, Caiming Xiong, Richard Socher, and Dragomir  
444 Radev. SummEval: Re-evaluating summarization evalu-  
445 ation. *Transactions of the Association for Computational*  
446 *Linguistics*, 9:391–409, 2021.
- 447 Suzanne Fricke. Semantic scholar. *Journal of the Medical*  
448 *Library Association: JMLA*, 106(1):145, 2018.
- 449 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-  
450 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu  
451 Chen. LoRA: Low-rank adaptation of large language mod-  
452 els. *arXiv preprint arXiv:2106.09685*, 2021.
- 453 Hang Jiang, Xiajie Zhang, Xubo Cao, Jad Kabbara, and Deb  
454 Roy. PersonaLLM: Investigating the ability of large lan-  
455 guage models to express personality traits. *arXiv preprint*  
456 *arXiv:2305.02547*, 2023.
- 457 Wojciech Kryściński, Bryan McCann, Caiming Xiong,  
458 and Richard Socher. Evaluating the factual consis-  
459 tency of abstractive text summarization. *arXiv preprint*  
460 *arXiv:1910.12840*, 2019.
- 461 Kuang-Huei Lee, Xinyun Chen, Hiroki Furuta, John Canny,  
462 and Ian Fischer. A human-inspired reading agent with  
463 gist memory of very long contexts. *arXiv preprint*  
464 *arXiv:2402.09727*, 2024.
- 465 Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen  
466 Xu, and Chenguang Zhu. G-Eval: NLG evaluation us-  
467 ing GPT-4 with better human alignment. *arXiv preprint*  
468 *arXiv:2303.16634*, 2023.
- 469 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape,  
470 Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost  
471 in the middle: How language models use long contexts.  
472 *Transactions of the Association for Computational Lin-*  
473 *guistics*, 12:157–173, 2024.
- OpenAI. GPT-4 technical report, 2023. 474
- Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo,  
475 Wangchunshu Zhou, Yuchen Eleanor Jiang, Chengfei  
476 Lv, and Huajun Chen. AutoAct: Automatic agent  
477 learning from scratch via self-planning. *arXiv preprint*  
478 *arXiv:2401.05268*, 2024. 479
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan  
480 Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill  
481 Qian, et al. ToolLLM: Facilitating large language mod-  
482 els to master 16000+ real-world APIs. *arXiv preprint*  
483 *arXiv:2307.16789*, 2023. 484
- Alireza Salemi, Sheshera Mysore, Michael Bendersky, and  
485 Hamed Zamani. LaMP: When large language models meet  
486 personalization. *arXiv preprint arXiv:2304.11406*, 2023. 487
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu.  
488 Character-LLM: A trainable agent for role-playing. *arXiv*  
489 *preprint arXiv:2310.10158*, 2023. 490
- Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper,  
491 Nicholas Lee, Shuo Yang, Christopher Chou, Banghua  
492 Zhu, Lianmin Zheng, Kurt Keutzer, et al. S-LoRA: Serv-  
493 ing thousands of concurrent LoRA adapters. *arXiv preprint*  
494 *arXiv:2311.03285*, 2023. 495
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales,  
496 David Dohan, Ed H Chi, Nathanael Schärli, and Denny  
497 Zhou. Large language models can be easily distracted by  
498 irrelevant context. In *International Conference on Machine*  
499 *Learning*, pages 31210–31227. PMLR, 2023. 500
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R  
501 Narasimhan, and Shunyu Yao. Reflexion: Language agents  
502 with verbal reinforcement learning. In *Thirty-seventh Con-*  
503 *ference on Neural Information Processing Systems*, 2023. 504
- Zhaoxuan Tan and Meng Jiang. User modeling in the era of  
505 large language models: Current research and future direc-  
506 tions. *arXiv preprint arXiv:2312.11518*, 2023. 507
- Xuewei Wang, Weiyan Shi, Richard Kim, Yoojung Oh, Sijia  
508 Yang, Jingwen Zhang, and Zhou Yu. Persuasion for good:  
509 Towards a personalized persuasive dialogue system for so-  
510 cial good. *arXiv preprint arXiv:1906.06725*, 2019. 511
- Yancheng Wang, Ziyang Jiang, Zheng Chen, Fan Yang,  
512 Yingxue Zhou, Eunah Cho, Xing Fan, Xiaojiang Huang,  
513 Yanbin Lu, and Yingzhen Yang. RecMind: Large language  
514 model powered agent for recommendation. *arXiv preprint*  
515 *arXiv:2308.14296*, 2023. 516
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Bar-  
517 ret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten  
518 Bosma, Denny Zhou, Donald Metzler, et al. Emer-  
519 gent abilities of large language models. *arXiv preprint*  
520 *arXiv:2206.07682*, 2022. 521
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma,  
522 Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-  
523 thought prompting elicits reasoning in large language mod-  
524 els. *Advances in neural information processing systems*,  
525 35:24824–24837, 2022. 526

527 Penghui Wei, Xuanhua Yang, Shaoguo Liu, Liang Wang,  
528 and Bo Zheng. Creator: Ctr-driven advertising text gen-  
529 eration with controlled pre-training and contrastive fine-  
530 tuning. *arXiv preprint arXiv:2205.08943*, 2022.

531 Weizhe Yuan, Graham Neubig, and Pengfei Liu. BARTScore:  
532 Evaluating generated text as text generation. *Advances in*  
533 *Neural Information Processing Systems*, 34:27263–27277,  
534 2021.

535 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-  
536 berger, and Yoav Artzi. BERTScore: Evaluating text gener-  
537 ation with BERT. *arXiv preprint arXiv:1904.09675*, 2019.

538 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris  
539 Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and  
540 Bill Dolan. DialoGPT: Large-scale generative pre-training  
541 for conversational response generation. *arXiv preprint*  
542 *arXiv:1911.00536*, 2019.

543 Kai Zhang, Fubang Zhao, Yangyang Kang, and Xiaozhong  
544 Liu. Memory-augmented LLM personalization with short-  
545 and long-term memory coordination. *arXiv preprint*  
546 *arXiv:2309.11696*, 2023.

547 Zhuosheng Zhang, Yao Yao, Aston Zhang, Xiangru Tang,  
548 Xinbei Ma, Zhiwei He, Yiming Wang, Mark Gerstein, Rui  
549 Wang, Gongshen Liu, et al. Igniting language intelligence:  
550 The hitchhiker’s guide from chain-of-thought reasoning to  
551 language agents. *arXiv preprint arXiv:2311.11797*, 2023.

552 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen,  
553 Heng-Tze Cheng, Ed H Chi, Quoc V Le, and Denny Zhou.  
554 Take a step back: Evoking reasoning via abstraction in  
555 large language models. *arXiv preprint arXiv:2310.06117*,  
556 2023.

557 Aakas Zhiyuli, Yanfang Chen, Xuan Zhang, and Xun Liang.  
558 BookGPT: A general framework for book recommenda-  
559 tion empowered by large language model. *arXiv preprint*  
560 *arXiv:2305.15673*, 2023.

561 Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jia-  
562 long Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing  
563 Chen, RuiPu Wu, Shuai Wang, et al. Agents: An open-  
564 source framework for autonomous language agents. *arXiv*  
565 *preprint arXiv:2309.07870*, 2023.



566 **A PSW Dataset Overview**

567 **Overview.** The PSW dataset is constructed using data from  
 568 the Semantic Scholar database [Fricke, 2018]. We first sel-  
 569 lected a subset of papers from Software Engineering pub-  
 570 lished after 2000, considering only papers with at least two  
 571 authors to ensure the feasibility of evaluating collaborative  
 572 writing scenarios. The collected papers were randomly split  
 573 into training, validation, and test subsets.<sup>4</sup> We performed the  
 574 split at the paper level to ensure that all tasks within the PSW  
 575 benchmark had consistent data splits. The summary of PSW  
 576 dataset statistics can be found in Table 5.

Statistic	Train	Validation	Test
# of Papers	1,744	500	500
# of Authors	6,461	1,655	1,280
Avg. Authors / Paper	4.05	3.16	3.25
Avg. History Papers / Author	63.47	75.34	92.21
Avg. Research Interests / Author	2.84	2.77	2.79
Avg. Title Length	97.03	95.54	96.16
Avg. Abstract Length	970.92	981.36	1,037.09
Avg. Research Question Length	470.57	398.22	442.31
Avg. References / Paper	60.24	54.85	58.93

Table 5: PSW Dataset Statistics with Train/Validation/Test Splits.

**B Comparison of Methods on PSW Dataset**

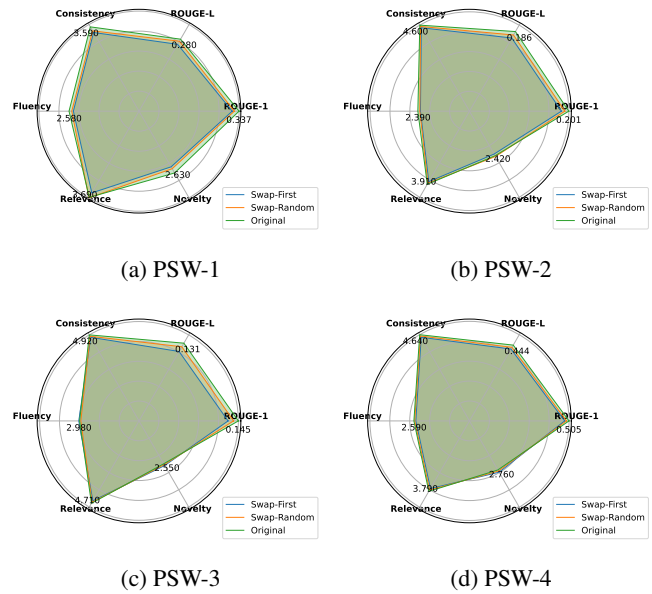


Figure 4: Impact of author order on the performance (ROUGE-1, ROUGE-L, Consistency, Fluency, Relevance, and Novelty) across three different models: **Original**, **Swap-Random**, and **Swap-First**. The **Original** model consistently achieves the highest scores across all datasets.

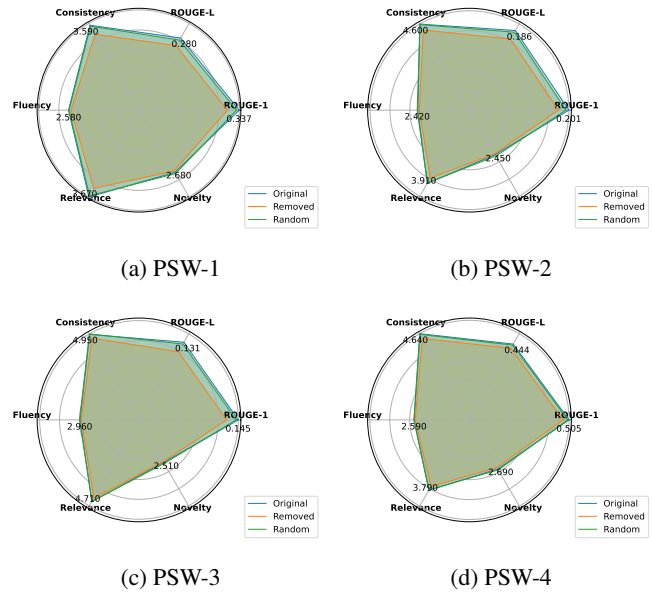


Figure 5: Impact of user profiling on the performance (ROUGE-1, ROUGE-L, Consistency, Fluency, Relevance, and Novelty) across three different models: **Original**, **Removed**, and **Random**. The **Original** model consistently achieves the highest scores across all datasets.

<sup>4</sup>We only used the test split in this paper since our method doesn't require model training.

## C Details of G-Eval Prompt

Task Description
You will be given one result generated for a science paper and several reference papers. Your task is to rate the result using the following criteria. Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.
Evaluation Criteria
<b>Consistency (1-5)</b> – the factual alignment between the result and the corresponding science paper. A factually consistent result contains only statements entailed by the source document.
<b>Fluency (1-3)</b> – the quality of the result in terms of grammar, spelling, punctuation, word choice, and sentence structure.
<b>Relevance (1-5)</b> – the selection of important content from the source. The result should include only important information from the source document.
<b>Novelty (1-3)</b> – the uniqueness and originality of the result in terms of concept, perspective, and creativity.
Evaluation Task
Now, you are working on evaluating this prediction: {Prediction Text}
Here are some ground truth results for comparison: [result <sub>1</sub> , result <sub>2</sub> , ...].
Instruction
Please evaluate the prediction using the above criteria.

Table 6: Prompt template for evaluating the G-Eval metric.

## D Prompts for LaMP Tasks

### D.1 Personalized Citation Identification (LaMP-1) Prompt

User Profile
Assuming you care a lot about these areas: <b>Keywords:</b> [keyword <sub>1</sub> , keyword <sub>2</sub> , keyword <sub>3</sub> , ...] <b>Topics:</b> [topics <sub>1</sub> , topics <sub>2</sub> , topics <sub>3</sub> , ...]
User History
I give you some titles of papers that you've written. Please imitate your reasons and recommend a paper citation for me. Each example consists of an abstract, the corresponding title, and a description of the writing style and keywords for that title.
Example 1
<b>Title:</b> {Title Text} <b>Abstract:</b> {Abstract Text} <b>Reason:</b> {Reason} <b>Citation:</b> [citation <sub>1</sub> , citation <sub>2</sub> , ...]
Example 2
<b>Title:</b> {Title Text} <b>Abstract:</b> {Abstract Text} <b>Reason:</b> {Reason} <b>Citation:</b> [citation <sub>1</sub> , citation <sub>2</sub> , ...] ...
Example k
<b>Title:</b> {Title Text} <b>Abstract:</b> {Abstract Text} <b>Reason:</b> {Reason} <b>Citation:</b> [citation <sub>1</sub> , citation <sub>2</sub> , ...]
Classification Task
Now you have written this title: <b>Title:</b> {Title Text}
Instruction
Please separately analyze the potential relevant connection of <b>Reference 1</b> and <b>Reference 2</b> to this title. You are citing from one of them. Please decide which one it would be: <b>Reference 1:</b> {option <sub>1</sub> } <b>Reference 2:</b> {option <sub>2</sub> } Just answer with [1] or [2] without explanation.

Table 7: Prompt template for the Personalized Citation Identification (LaMP-1) task.

## D.2 Personalized News Categorization (LaMP-2) Prompt

<b>User Profile</b>
Assuming you care a lot about these areas: <b>Keywords:</b> [keyword <sub>1</sub> , keyword <sub>2</sub> , keyword <sub>3</sub> , ...] <b>Topics:</b> [topics <sub>1</sub> , topics <sub>2</sub> , topics <sub>3</sub> , ...]
<b>User History</b>
I give you some titles and articles that you've written with category. Please imitate your reasons for giving this category. Each example consists of an abstract, the corresponding title, and a category of it.
<b>Example 1</b>
<b>Article:</b> {Article Text} <b>Title:</b> {Title Text} <b>Reason:</b> {Reason} <b>Category:</b> [category <sub>1</sub> , category <sub>2</sub> , ...]
<b>Example 2</b>
<b>Article:</b> {Article Text} <b>Title:</b> {Title Text} <b>Reason:</b> {Reason} <b>Category:</b> [category <sub>1</sub> , category <sub>2</sub> , ...] ...
<b>Example k</b>
<b>Article:</b> {Article Text} <b>Title:</b> {Title Text} <b>Reason:</b> {Reason} <b>Category:</b> [category <sub>1</sub> , category <sub>2</sub> , ...]
<b>Classification Task</b>
Now you have written this article with the title: <b>Article:</b> {Article Text} <b>Title:</b> {Title Text}
<b>Instruction</b>
Which category does this article relate to among the following categories? <b>Category 1:</b> {option <sub>1</sub> } <b>Category 2:</b> {option <sub>2</sub> } ... <b>Category K:</b> {option <sub>N</sub> } Just answer with the category name without further explanation.

Table 8: Prompt template for the Personalized News Categorization (LaMP-2) task.

## D.3 Personalized Product Rating (LaMP-3) Prompt

<b>User Profile</b>
Assuming you have written product reviews with the following characteristics: <b>Most Common Rating:</b> {score <sub>most</sub> } <b>Rating Patterns:</b> [pattern <sub>1</sub> , pattern <sub>2</sub> , ...]
<b>User History</b>
I provide you with some product reviews you've written, along with their corresponding ratings. Please imitate your reasoning for assigning these ratings. Each example consists of a product review and its rating.
<b>Example 1</b>
<b>Product Review:</b> {Review Text} <b>Rating:</b> {Rating}
<b>Example 2</b>
<b>Product Review:</b> {Review Text} <b>Rating:</b> {Rating} ...
<b>Example k</b>
<b>Product Review:</b> {Review Text} <b>Rating:</b> {Rating}
<b>Rating Task</b>
Now you have written this new product review: <b>Product Review:</b> {Review Text} Based on the review, please analyze its sentiment and how much you like the product.
<b>Instruction</b>
Follow your previous rating habits and these instructions: <ul style="list-style-type: none"> <li>• If you feel satisfied with this product or have concerns but it's good overall, it should be rated 5.</li> <li>• If you feel good about this product but notice some issues, it should be rated as 4.</li> <li>• If you feel OK but have concerns, it should be rated as 3.</li> <li>• If you feel unsatisfied with this product but it's acceptable for some reason, it should be rated as 2.</li> <li>• If you feel completely disappointed or upset, it should be rated 1.</li> </ul>
Your most common rating is {score <sub>most</sub> }. You must follow this rating pattern faithfully and answer with the rating without further explanation.

Table 9: Prompt template for the Personalized Product Review Rating (LaMP-3) task.

#### D.4 Personalized News Headline Generation (LaMP-4) Prompt

<b>User Profile</b>
Assuming you have written headlines with the following characteristics:
<b>Writing Style:</b> [style <sub>1</sub> , style <sub>2</sub> , ...]
<b>Content Patterns:</b> [patterns <sub>1</sub> , patterns <sub>2</sub> , ...]
<b>User History</b>
I will provide you with some news articles along with the headlines you've written for them. Please imitate your writing style and content patterns when generating a new headline. Each example consists of a news article and its corresponding headline.
<b>Example 1</b>
<b>Article:</b> {Article Text}
<b>Headline:</b> {Headline}
<b>Example 2</b>
<b>Article:</b> {Article Text}
<b>Headline:</b> {Headline}
...
<b>Example k</b>
<b>Article:</b> {Article Text}
<b>Headline:</b> {Headline}
<b>Generation Task</b>
Now that you have been given this news article:
<b>Article:</b> {Article Text}
<b>Instruction</b>
Please write a headline following your previous writing styles and habits. If you have written headlines with similar content, you could reuse those headlines and mimic their content.

Table 10: Prompt template for the Personalized News Headline Generation (LaMP-4) task.

#### D.5 Personalized Scholarly Title Generation (LaMP-5) Prompt

<b>User Profile</b>
Assuming you have written scholarly titles with the following characteristics:
<b>Writing Style:</b> [style <sub>1</sub> , style <sub>2</sub> , ...]
<b>Title Patterns:</b> [pattern <sub>1</sub> , pattern <sub>2</sub> , ...]
<b>User History</b>
I will provide you with some research paper abstracts along with the titles you've written for them. Please imitate your writing style and title patterns when generating a new title. Each example consists of a paper abstract and its corresponding title.
<b>Example 1</b>
<b>Abstract:</b> {Abstract Text}
<b>Title:</b> {Title}
<b>Example 2</b>
<b>Abstract:</b> {Abstract Text}
<b>Title:</b> {Title}
...
<b>Example k</b>
<b>Abstract:</b> {Abstract Text}
<b>Title:</b> {Title}
<b>Generation Task</b>
Now that you have been given this paper abstract:
<b>Abstract:</b> {Abstract Text}
<b>Instruction</b>
Please write a title following your previous style and habits, keeping it clear, accurate, and concise.

Table 11: Prompt template for the Personalized Scholarly Title Generation (LaMP-5) task.

## D.6 Personalized Email Subject Generation (LaMP-6) Prompt

<b>User Profile</b>
Assuming you care a lot about these areas: <b>Keywords:</b> [keyword <sub>1</sub> , keyword <sub>2</sub> , keyword <sub>3</sub> , ...] <b>Topics:</b> [topics <sub>1</sub> , topics <sub>2</sub> , topics <sub>3</sub> , ...]
<b>User History</b>
Let's say there are some emails you've written. Please mimic the style of these examples. Each example consists of email content, the corresponding subject, and a description of the writing style for that title.
<b>Example 1</b>
<b>Content:</b> {Email Content} <b>Writing Style:</b> {Style} <b>Subject:</b> {Email Subject}
<b>Example 2</b>
<b>Content:</b> {Email Content} <b>Writing Style:</b> {Style} <b>Subject:</b> {Email Subject}
...
<b>Example k</b>
<b>Content:</b> {Email Content} <b>Writing Style:</b> {Style} <b>Subject:</b> {Email Subject}
<b>Generation Task</b>
Now that you have been given this email content: <b>Content:</b> {Email Content}
<b>Instruction</b>
Write a title following your previous style and habits. Just answer with the subject without further explanation.

Table 12: Prompt template for the Personalized Email Subject Generation (LaMP-6) task.

## D.7 Personalized Tweet Paraphrasing (LaMP-7) Prompt

<b>User Profile</b>
Assuming you have written tweets with the following characteristics: <b>Writing Style:</b> [style <sub>1</sub> , style <sub>2</sub> , ...] <b>Tone:</b> [tone <sub>1</sub> , tone <sub>2</sub> , ...] <b>Length:</b> [length <sub>1</sub> , length <sub>2</sub> , ...]
<b>User History</b>
I will provide you with some original tweets along with the paraphrased versions you've written for them. When paraphrasing a new tweet, please imitate your writing style, tone, and typical length. Each example consists of an original tweet and its paraphrased version.
<b>Example 1</b>
<b>Original Tweet:</b> {Tweet Text} <b>Paraphrased Tweet:</b> {Paraphrased Text}
<b>Example 2</b>
<b>Original Tweet:</b> {Tweet Text} <b>Paraphrased Tweet:</b> {Paraphrased Text}
...
<b>Example k</b>
<b>Original Tweet:</b> {Tweet Text} <b>Paraphrased Tweet:</b> {Paraphrased Text}
<b>Generation Task</b>
Now that you have been given this tweet: <b>Original Tweet:</b> {Tweet Text}
<b>Instruction</b>
Please paraphrase it with the following instructions:
<ul style="list-style-type: none"> <li>You must use tweet styles and tones.</li> <li>You must keep it faithful to the given tweet with similar keywords and length.</li> </ul>

Table 13: Prompt template for the Personalized Tweet Paraphrasing (LaMP-7) task.

594 **E Prompts for PSW Tasks**

595 **E.1 Research Interests Generation (UP-0) Prompt**

<b>User History</b>
I will provide you with some research papers you've authored. Please summarize your top research interests based on these papers. Each paper consists of a title and abstract.
<b>Paper 1</b>
<b>Title:</b> {Title Text}
<b>Abstract:</b> {Abstract Text}
<b>Paper 2</b>
<b>Title:</b> {Title Text}
<b>Abstract:</b> {Abstract Text}
...
<b>Paper k</b>
<b>Title:</b> {Title Text}
<b>Abstract:</b> {Abstract Text}
<b>Instruction</b>
Please summarize your top three research interests based on the provided papers in the following format: <b>Research Interests:</b> [interest <sub>1</sub> , interest <sub>2</sub> , interest <sub>3</sub> , ...]

Table 14: Prompt template for the Research Interests Generation (UP-0) task.

596 **E.2 Personalized Research Paper Title Generation (PSW-1) Prompt**

596  
597

<b>User Profile</b>
Assuming you are an expert researcher with the following research interests: <b>Research Interests:</b> [interest <sub>1</sub> , interest <sub>2</sub> , interest <sub>3</sub> , ...]
<b>User History</b>
Here are some titles and abstracts from papers you have authored:
<b>Paper 1</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Paper 2</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
...
<b>Paper k</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Brainstorm Task</b>
Here are some related papers for reference, each with a title: <b>Reference 1:</b> {Title} <b>Reference 2:</b> {Title}
...
<b>Reference N:</b> {Title}
<b>Instruction</b>
Considering your research interests, previous works, and reference papers, please brainstorm the most promising title for your new research paper.

Table 15: Prompt template for the Personalized Research Paper Title Generation (PSW-1) task.

### E.3 Research Question Generation (PSW-2) Prompt

<b>User Profile</b>
Assuming you are an expert researcher with the following research interests:
<b>Research Interests:</b> [interest <sub>1</sub> , interest <sub>2</sub> , interest <sub>3</sub> , ...]
<b>User History</b>
Here are some titles and abstracts from papers you have authored:
<b>Paper 1</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Paper 2</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
...
<b>Paper k</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Brainstorm Task</b>
Now you are working on a new paper with the following title:
<b>Title:</b> {Title}
<b>Instruction</b>
Considering the title and research background, please propose the top 3 research questions you aim to address in this new paper.

Table 16: Prompt template for the Research Question Generation (PSW-2) task.

### E.4 Paper Abstract Generation (PSW-3) Prompt

<b>User Profile</b>
Assuming you are an expert researcher with the following research interests:
<b>Research Interests:</b> [interest <sub>1</sub> , interest <sub>2</sub> , interest <sub>3</sub> , ...]
<b>User History</b>
Here are some titles and abstracts from papers you have authored:
<b>Paper 1</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Paper 2</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
...
<b>Paper k</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Generation Task</b>
Now you are working on a new paper with the following title:
<b>Title:</b> {Title}
And you are focusing on solving the following research questions: [question <sub>1</sub> , question <sub>2</sub> , ...]
<b>Instruction</b>
Considering the title, research questions, and your writing style in previous abstracts, please write an abstract for this new paper.

Table 17: Prompt template for the Paper Abstract Generation (PSW-3) task.

<b>User Profile</b>
Assuming you are an expert researcher with the following research interests: <b>Research Interests:</b> [interest <sub>1</sub> , interest <sub>2</sub> , interest <sub>3</sub> , ...]
<b>User History</b>
Here are some titles and abstracts from papers you have authored:
<b>Paper 1</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Paper 2</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
...
<b>Paper k</b>
<b>Title:</b> {Title}
<b>Abstract:</b> {Abstract}
<b>Generation Task</b>
Now, you are working on a new paper with the following abstract: <b>Abstract:</b> {Abstract}
And you are focusing on solving the following research questions: [question <sub>1</sub> , question <sub>2</sub> , ...]
<b>Instruction</b>
Considering the abstract and your title writing style in previous papers, please generate a title for this new paper. The title should be clear and concise and reflect the main topic of the abstract as well as your research questions.

Table 18: Prompt template for the Paper Title Generation (PSW-4) task.