

LaMP: When Large Language Models Meet Personalization

Alireza Salemi¹, Sheshera Mysore¹, Michael Bendersky², Hamed Zamani¹

¹University of Massachusetts Amherst

²Google Research

{asalemi, smysore, zamani}@cs.umass.edu

bemike@google.com

The LaMP Benchmark: <http://lamp-benchmark.github.io/>

Abstract

This paper highlights the importance of personalization in the current state of natural language understanding and generation and introduces the LaMP benchmark — a novel benchmark for training and evaluating language models for producing personalized outputs. LaMP offers a comprehensive evaluation framework with diverse language tasks and multiple entries for each user profile. It consists of seven personalized tasks, spanning three classification and four text generation tasks. We also propose a retrieval augmentation approach that retrieves personalized items from user profiles to construct personalized prompts for large language models. Our baseline zero-shot and fine-tuned model results indicate that LMs utilizing profile augmentation outperform their counterparts that do not factor in profile information.

1 Introduction

As natural language processing (NLP) systems evolve, personalization has emerged as a key factor in meeting the user’s expectations for tailored experiences that align with their unique needs and preferences. While personalization has been widely studied by various communities, including the information retrieval (IR) and human-computer interaction (HCI) communities, often with applications to search engines and recommender systems (Fowler et al., 2015; Xue et al., 2009; Naumov et al., 2019) – its exploration in NLP has been limited. However, the importance of personalization for text classification and generation tasks has been highlighted in recent work of Flek (2020) and Dudy et al. (2021). This work also notes the potential of personalization for centering users and creating accessible and inclusive systems. This optimism has also been reflected in the recent UserNLP’22 workshop (Huang et al., 2022).

In tandem, the recent introduction of large language models (LLMs), such as GPT4 (OpenAI,

2023), has revolutionized natural language processing. With recent work has also highlighted the benefits (and harms) of personalizing LLMs (Kirk et al., 2023). Despite this and the importance of personalization in many real-world problems, developing and evaluating LLMs for producing personalized responses remain understudied. Therefore, in this paper, we underscore the importance of personalization in shaping the future of NLP and take the first step towards developing and evaluating personalization in the context of large language models by proposing the LaMP benchmark.¹

While many existing well-known NLP benchmarks, such as GLUE (Wang et al., 2018), SuperGLUE (Wang et al., 2019), KILT (Petroni et al., 2021), and GEM (Gehrmann et al., 2021) have significantly progressed the NLP frontier, they have often taken the dominant NLP approach of “one-size-fits-all” to modeling and evaluation, and do not allow the development of models that adapt to the specific needs of end users – limiting extensive research on personalization in NLP tasks. In contrast, LaMP offers a comprehensive evaluation framework incorporating diverse language tasks that require personalization. LaMP consists of three personalized text classification tasks: (1) Personalized Citation Identification (binary classification), (2) Personalized News Categorization (categorical classification with 15 categories), and (3) Personalized Product Rating (ordinal classification from 1 to 5-star rating for e-commerce products). Moreover, LaMP includes four text generation datasets: (4) Personalized News Headline Generation, (5) Personalized Scholarly Title Generation, (6) Personalized Email Subject Generation, and (7) Personalized Tweet Paraphrasing. For each of these seven tasks, we utilize two different data splitting settings: (a) user-based data splitting, in which each user appear only in one of the train, validation, or test sets. This setting enables us to study

¹LaMP stands for Language Model Personalization.

personalization for unseen users. And (b) time-based data splitting, in which the users are shared between the aforementioned sets and the data is split based on time. This setting enables us to study personalization for future interactions of known users. Therefore, LaMP provides a rich environment for developing personalized NLP models.

For personalizing the language model outputs, a straightforward solution is to incorporate the user profile into a language model prompt. However, user profiles are often large and exceed the length limitations of large language models. Even if we relax such limitations as the technology evolves, the cost of processing large input sequences is considerable. Therefore, we propose a personalized retrieval augmentation solution, where for each test input, we retrieve personalized items from the user profile to be included in the LM prompt. We demonstrate that using this approach, the performance of language models improves on all datasets in the LaMP benchmark. Based on this retrieval augmentation solution, we evaluate different retrievers for personalized prompt construction and establish benchmark results for fine-tuned and zero-shot language models. The LaMP benchmark, the data construction and evaluation scripts, and leaderboard are publicly available for research purposes: <http://lamp-benchmark.github.io/>.

2 Personalizing LLM Outputs

2.1 Problem Formulation

The LaMP benchmark considers every data sample as an individual user, with a collection of user records associated with each sample for all tasks. These user records can facilitate the personalization of language models based on user-specific data. Consequently, each data sample can be partitioned into three distinguishable components: an input sequence that serves as the model’s input, a target output that the model is expected to produce, and a profile that encapsulates any auxiliary information that can be employed to personalize the model according to the user’s specific preferences or requirements. Therefore, all tasks within the LaMP benchmark are formalized as follows: for a given textual input x , the goal is to develop a model M that generates personalized output y for user u . One can model this task as $\arg \max_y p(y|x, u)$. For each user u , the model M can take advantage of $P_u = \{(x_{u1}, y_{u1}), (x_{u2}, y_{u2}), \dots, (x_{um_u}, y_{um_u})\}$ where each (x_{ui}, y_{ui}) denotes a pair of input and

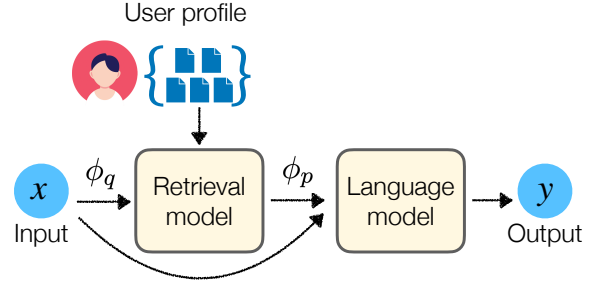


Figure 1: An overview of the retrieval-augmented method for personalizing LLMs. ϕ_q and ϕ_p represent query and prompt construction functions.

personalized output for user u in the same format as x and y . Each task is accompanied by a training set that does not share any user with the test set.

2.2 A Retrieval Augmentation Approach for Personalizing LLMs

In the tasks at hand, each user profile consists of a (potentially large) collection of data points pertaining to the user. Given the inherent context length constraint of language models in addition to their efficiency and cost, it is only practical to incorporate a subset of these data points as input prompts. Moreover, it is important to note that not all entries within a user profile are necessarily relevant to the specific task the user aims to accomplish. Consequently, we propose the development of solutions through retrieval augmentation. This framework selectively extracts pertinent information from the user profile that are relevant to the current unseen test case. An overview of the method is shown in Figure 1.

To achieve personalization for a given sample (x_i, y_i) associated with user u , we employ three primary components: (1) a query generation function ϕ_q that transforms the input x_i into a query q for retrieving from the user u ’s profile, (2) a retrieval model $\mathcal{R}(q, P_u, k)$ that accepts a query q , a user profile P_u and retrieves k most pertinent entries from the user profile, and (3) a prompt construction function ϕ_p that assembles a prompt for user u based on input x_i and the retrieved entries. Consequently, the input \bar{x}_i for the language model is derived using the following formulation:

$$\bar{x}_i = \phi_p(x_i, \mathcal{R}(\phi_q(x_i), P_u, k)) \quad (1)$$

where we use (\bar{x}_i, y_i) to train or evaluate the language models. In this paper, we use Contriever (Izcard et al., 2022), BM25 (Robertson et al.,

1995), Recency in which we select the most recent items from the user profile (only for the time-based separation setting), and a random selection from user’s profile to model \mathcal{R} . The implementation of the ϕ_p function is depicted in Table 4 in Appendix D. For the ϕ_q function, we concatenate the non-template parts of each input, which can be seen in Figure 3 in Appendix C to obtain the query.

3 The LaMP Benchmark

The LaMP benchmark aim at assessing the efficacy of language models in producing personalized outputs based on user-specific information for seven diverse tasks. These tasks can be categorized as either personalized text classification or personalized text generation tasks:

- **Personalized Text Classification**

- (1) Personalized Citation Identification
- (2) Personalized News Categorization
- (3) Personalized Product Rating

- **Personalized Text Generation**

- (4) Personalized News Headline Generation
- (5) Personalized Scholarly Title Generation
- (6) Personalized Email Subject Generation
- (7) Personalized Tweet Paraphrasing

This benchmark creates each dataset in two different settings: 1) user-based separation and 2) time-based separation.

User-based Separation. In the context of user-based separation, individuals are partitioned into training, validation, and test sets. Notably, the temporal aspect is disregarded in this scenario, whereby the selection of profile items, inputs, and outputs is randomized rather than time-dependent.

Time-based Separation. In the time-based separation setting, the distribution of users spans across the train, validation, and test sets. Moreover, the selection process for profile items, inputs, and outputs is influenced by temporal factors. Specifically, the most recent user information is chosen to construct the input and output for train, validation, and test sets, while older information is utilized as profile items.

In the following sections, we introduce our data collection approach for each of these tasks and each setting.

3.1 Task 1: Personalized Citation Identification

This task assesses the capacity of a language model to establish a connection between two papers’ titles based on the user profile. In more detail, if the user u writes a paper on a topic x , the language model should determine what papers they cite. Figure 3 in Appendix C illustrates an instance of this task, wherein the model is required to predict the association between a paper’s title written by u and two candidate papers, which may serve as references. Therefore, this is a binary classification task.

Data Collection. To generate data samples, we leverage the Citation Network Dataset (V14) (Tang et al., 2008), which comprises information on scientific papers, authors, and citations. We select all papers from this dataset that meet the following criteria: 1) they are written in English, 2) they contain at least one reference and one author, and 3) they include an abstract. Subsequently, we group papers based on their authors and only consider authors who have written at least 50 papers. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

Profile Specification. For this task, the profile of each user encompasses all the papers they have authored. We retain only the title and abstract of each paper in the user’s profile. We exclude the paper selected to generate the input sequence for this task from the user’s profile.

Evaluation. To evaluate performance on this task, we partition the data into train, validation, and test sets. A summary of the dataset statistics is provided in Table 9 in Appendix F. As this is a balanced binary classification task, we adopt accuracy as the sole evaluation metric.

3.2 Task 2: Personalized News Categorization

This task aims to assess the capability of a language model to classify news articles written by a user (journalist) u . An illustration of this task is presented in Figure 3 in Appendix C. Given a news article x written by u , the language model must predict its category from the set of available categories based on the user’s past articles.

Data Collection. To construct our dataset for this task, we leverage the news categorization dataset

(Misra, 2022; Misra and Grover, 2021) obtained from the HuffPost website². However, we filter out some categories and merge similar ones to form a more concise set of categories, as described in Appendix B. Next, we group articles by their author, taking only the first author in cases where there are multiple authors, and retain only authors with a minimum of specific number articles who have published in at least three different categories. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

Profile Specification. For this task, the user profile consists of the articles written by the user along with their respective categories. However, we exclude the article that is selected as the input for the task from the user profile.

Evaluation. For this task, we create three sets: training, validation, and testing, and their statistics are reported in Table 9 in Appendix F. Given that this task involves multi-class classification, we employ accuracy and macro-averaged F1-score as the evaluation metrics to measure the model’s performance.

3.3 Task 3: Personalized Product Rating

This task evaluates the language model’s ability to predict the rating that a user u has given to a product based on the review written by u for the product. The user’s profile, as illustrated in Figure 3 in Appendix C, serves as a basis for predicting an ordinal score with a range from 1 to 5, and exclusively consisting of integer values.

Data Collection. In this task, we create our dataset by leveraging the Amazon Reviews Dataset (Ni et al., 2019). We filtered out users (i.e., amazon customers who have written reviews) who have written less than 100 and the 1% users with the most reviews as outliers. Since the Amazon Reviews dataset is quite extensive, we randomly sampled a subset of users from the dataset. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

Profile Specification. In this task, the profile refers to a user’s other reviews and their respective assigned ratings. It is worth noting that the article selected as the input for the task does not

contribute to the user’s profile. Instead, the profile is solely derived from the remaining reviews authored by the user.

Evaluation. To evaluate the performance of our model, we randomly split our dataset into training, validation, and test sets. The relevant statistics of our dataset are presented in Table 9 in Appendix F. Given that our task is an ordinal multi-class classification problem, we employ RMSE and MAE as the primary evaluation metrics.

3.4 Task 4: Personalized News Headline Generation

This task aims to evaluate the language model’s capability to generate a headline for a news article written by a user u . An example of this task is presented in Figure 3 in Appendix C. Specifically, the task involves providing the language model with a news article and requesting it to generate a headline that accurately reflects the user’s interests and writing style, as captured in their profile. This task assesses the model’s ability to produce headlines that are informative and personalized, based on the user’s profile.

Data Collection. To construct our dataset for this task, we leverage the News Categorization dataset (Misra, 2022; Misra and Grover, 2021) from the HuffPost website³. The dataset provides author information for each article and is used to group articles by their respective authors. We use the same method for filtering out authors as Section 3.2. In cases where an article has multiple authors, we assign it only to the first author. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

Profile Specification. For this task, we define the user profile as the collection of previous articles and their corresponding headlines written by the same author. We exclude the selected article for the sample input from the user’s profile.

Evaluation. To evaluate the performance of the model on the headline generation task, we create a train, validation, and test set. The dataset statistics are presented in Table 9 in Appendix F. To measure the quality of the generated headlines, we adopt Rouge-1 and Rouge-L as evaluation metrics, which have been commonly used in previous work for

²<https://www.huffpost.com/>

³<https://www.huffpost.com/>

headline generation tasks (Yamada et al., 2021; Panthaplackel et al., 2022). These metrics capture the overlap between the generated headline and the ground-truth reference headlines based on n-gram overlap and longest common subsequence (LCS) matching.

3.5 Task 5: Personalized Scholarly Title Generation

This task evaluates the language model’s capability to generate a title for a research paper, taking into account the other papers authored by the user. An example of this task is illustrated in Figure 3 in Appendix C, where the language model is presented with an abstract of a paper and is required to generate a title that aligns with the user’s profile.

Data Collection. Similar to Section 3.1, we leverage the Citation Network Dataset (V14) (Tang et al., 2008) that includes information about scientific papers, authors, and citations to construct our dataset. We only kept the papers that meet the following criteria: 1) written in English, 2) have at least one reference and one author, and 3) have an abstract. Then, we group papers by their authors and only consider authors who have published at least 50 papers. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

Profile Specification. In this task, the user profile comprises all the papers authored by the user. We extract only the title and abstract of each paper to form the user’s profile. Notably, we exclude the paper chosen to create the task input from the user profile.

Evaluation. For this task, we generate a training, validation, and test dataset. We report the statistics of this dataset in Table 9 in Appendix F. Similar to the news headline generation task, we adopt Rouge-1 and Rouge-L (Lin, 2004) as evaluation metrics.

3.6 Task 6: Personalized Email Subject Generation

The primary objective of this task is to evaluate the language model’s proficiency in generating an appropriate email subject based on the user’s writing style. Figure 3 in Appendix C provides an example of the task, which involves providing an email as input to the language model and requesting it to generate a corresponding subject that accurately

reflects the content of the email while aligning with the user’s writing style.

Data Collection. In this study, we adopt the Avocado Research Email Collection (Oard, Douglas et al., 2015) as the primary dataset for our task. To curate the dataset, we first perform a filtering step where we exclude emails with subject lengths of fewer than five words and content lengths of fewer than 30 words. Next, we group the emails based on their sender’s email address, retaining only those from users with email frequencies ranging between 10 to 200 emails. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

Profile Specification. For this task, we define a user’s profile as the entire set of emails authored by that user. To construct the profile, we extract both the subject and content of each email, excluding the email selected to serve as the input for the given task.

Evaluation. In this study, we create distinct training, validation, and test datasets to facilitate model development and evaluation. We provide an overview of these datasets’ key statistics in Table 9 in Appendix F. We adopt the Rouge-1 and Rouge-L metrics (Lin, 2004) for evaluation.

3.7 Task 7: Personalized Tweet Paraphrasing

The proposed task aims to evaluate the language model’s capacity to generate a paraphrased version of a tweet, considering the writing style of the user. An illustration of the task is presented in Figure 3 in Appendix C. Specifically, the language model is presented with a tweet and instructed to generate a corresponding paraphrase that captures the essence of the original tweet while aligning with the writing style of the user.

Data Collection. In this task, we utilize the Sentiment140 dataset (Go et al., 2009) as our tweet collection set. To ensure that the collected tweets are of adequate length, we only retain tweets containing at least 10 words. We then group the tweets based on the user ID and filter out users with fewer than 10 tweets. Additionally, we use ChatGPT (i.e., gpt3.5-turbo)⁴ to create inputs for this task. The details of data creation approach for this dataset in both user- and time-based separation are reported in Appendix A.

⁴<https://openai.com/blog/chatgpt>

Profile Specification. We construct user profiles using all the tweets that a user has posted, excluding the tweet that was selected to form the input to the task. To this end, we only retain the tweet text and disregard other metadata associated with each tweet.

Evaluation. We partition the collected dataset into three distinct subsets: train, validation, and test. We report the key statistics of this dataset in Table 9 in Appendix F. In line with prior research on paraphrase generation (Zhou and Bhat, 2021), we utilize Rouge-1 and Rouge-L (Lin, 2004) as the primary evaluation metrics for this task.

4 Experiments

This section describes our experiments on LaMP and provides benchmark results.

4.1 Experimental Setup

To train our models, we leverage the AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of 5×10^{-5} . We set 5% of the total training steps as warmup steps using a linear warmup scheduler. We also incorporate a weight decay of 10^{-4} to prevent overfitting during training. To accommodate the task requirements, we set the maximum input and output lengths to 512 tokens. We train our generation and classification models for 20 and 10 epochs, respectively. We utilize a FlanT5-base (Chung et al., 2022) model for all experiments, unless explicitly stated otherwise. We employ beam search (Freitag and Al-Onaizan, 2017) with a beam size of 4 in all experiments to improve the model’s ability to generate high-quality predictions.

4.2 Personalized Classification and Generation Results for Fine-Tuned Language Models

In this section, we establish baseline results for fine-tuned language models. We also investigate the impact of employing various retrieval techniques and the effect of retrieving different quantities of entries from a user’s profile. This analysis aims to provide insights into the efficacy of diverse retrieval methods and the potential benefits of adjusting the number of retrieved entries for personalization tasks.

The Impact of Retrievers on the End-to-End Performance. In our experiments, we employ four retrieval approaches for implementing \mathcal{R} : 1) a baseline random selector from the user profile, 2)

BM25 (Robertson et al., 1995), 3) Contriever (Izacard et al., 2022), and 4) Recency, in which we select the latest item in the user profile based on time (only for time-based separation setting). BM25 is considered as a robust and strong term-matching retrieval model and Contriever is a pretrained dense retrieval model. FlanT5-base is fine-tuned in all the experiments in this section as the language model that generates personalized output. The results of this experiment are shown in Table 1 for user-based separation and Table 2 for time-based separation. The results suggest that personalization improves the performance for all the text classification and generation tasks within the LaMP benchmark. In majority of cases, even a random selection of documents from the user profile leads to performance improvements.

When retrieving one document per user for personalizing the language model’s output, Contriever demonstrates the best performance in almost all classification datasets in user-based separation setting (i.e., LaMP-1U, LaMP-2U, LaMP-3U, LaMP-1T, and LaMP-2T). Recency only outperforms Contriever in the LaMP-3T classification task, which is reasonable because other studies have shown recency’s importance in product suggestion systems (Fader et al., 2005; Reinartz and Kumar, 2000, 2003). For text generation, Contriever performs best for Personalization News Headline Generation (LaMP-4U) and Personalized Tweet Paraphrasing (LaMP-7U) in user-based separation setting. For Email Generation and Scholarly Title Generation tasks (LaMP-5U and LaMP-6U), BM25 demonstrates superior performance. Both BM25 and Contriever outperform a random profile selector in all LaMP datasets. For the time-based separation setting, Contriever outperforms other methods in all generation tasks except News Headline Generation (LaMP-4T), where recency performs better.

Generally, the results we achieved indicates that merely incorporating information from the user profile into the input is not sufficient, but rather selecting the most relevant or recent information is crucial. This underscores the importance of careful consideration in selecting and incorporating pertinent user profile elements in language model prompts.

The Impact of k (the Number of Retrieved Items from each User Profile) on the End-to-End Performance. Each sample within this benchmark consists of a substantial number of user

Dataset	Metric	FlanT5-base (fine-tuned)				
		Non-Personalized	Untuned profile, $k = 1$			Tuned profile
			Random	BM25	Contriever	
LaMP-1U: Personalized Citation Identification	Accuracy	0.518	0.598	0.649	0.688	0.734
LaMP-2U: Personalized News Categorization	Accuracy	0.674	0.699	0.718	0.729	0.763
	F1	0.499	0.522	0.546	0.555	0.614
LaMP-3U: Personalized Product Rating	MAE	0.275	0.284	0.258	0.248	0.246
	RMSE	0.581	0.602	0.573	0.563	0.565
LaMP-4U: Personalized News Headline Generation	ROUGE-1	0.153	0.162	0.167	0.173	0.186
	ROUGE-L	0.140	0.148	0.153	0.159	0.171
LaMP-5U: Personalized Scholarly Title Generation	ROUGE-1	0.418	0.409	0.440	0.431	0.450
	ROUGE-L	0.378	0.371	0.399	0.393	0.409
LaMP-6U: Personalized Email Subject Generation	ROUGE-1	0.379	0.486	0.586	0.572	0.587
	ROUGE-L	0.358	0.470	0.570	0.558	0.575
LaMP-7U: Personalized Tweet Paraphrasing	ROUGE-1	0.509	0.514	0.521	0.524	0.528
	ROUGE-L	0.455	0.460	0.468	0.471	0.475

Table 1: The personalized text classification and generation results for a fine-tuned language model (i.e., FlanT5-base) on the test set of user-based separation setting. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3U dataset. k denotes the number of documents retrieved for personalizing language model outputs. For more information about tuned profile, you can check the Table 5 in Appendix E.

Dataset	Metric	FlanT5-base (fine-tuned)					Tuned profile
		Non-Personalized	Untuned profile, $k = 1$				
			Random	BM25	Contriever	Recency	
LaMP-1T: Personalized Citation Identification	Accuracy	0.628	0.657	0.682	0.688	0.691	0.714
LaMP-2T: Personalized News Categorization	Accuracy	0.762	0.794	0.783	0.815	0.800	0.806
	F1	0.574	0.634	0.613	0.656	0.645	0.659
LaMP-3T: Personalized Product Rating	MAE	0.280	0.279	0.278	0.281	0.279	0.266
	RMSE	0.615	0.612	0.614	0.606	0.608	0.598
LaMP-4T: Personalized News Headline Generation	ROUGE-1	0.159	0.169	0.171	0.176	0.173	0.177
	ROUGE-L	0.145	0.155	0.157	0.162	0.158	0.162
LaMP-5T: Personalized Scholarly Title Generation	ROUGE-1	0.462	0.460	0.471	0.472	0.466	0.479
	ROUGE-L	0.416	0.414	0.423	0.426	0.420	0.431
LaMP-6T: Personalized Email Subject Generation	ROUGE-1	0.479	0.525	0.537	0.545	0.532	0.547
	ROUGE-L	0.463	0.507	0.522	0.530	0.518	0.533
LaMP-7T: Personalized Tweet Paraphrasing	ROUGE-1	0.462	0.505	0.508	0.505	0.503	0.516
	ROUGE-L	0.416	0.456	0.457	0.455	0.453	0.465

Table 2: The personalized text classification and generation results for a fine-tuned language model (i.e., FlanT5-base) on the test set of time-based separation setting. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3U dataset. k denotes the number of documents retrieved for personalizing language model outputs. For more information about tuned profile, you can check the Table 7 in Appendix E.

profile entries. As such, exploring the impact of incorporating multiple entries to augment the input of the language model can provide valuable insights into addressing the unresolved challenges

posed by this benchmark. Figure 2 depicts the outcome of applying the optimal retriever from Tables 1 and 2 across various tasks, while varying the number of retrieved entries from user profiles.

Dataset	Metric	User-based Separation				Time-based Separation			
		Non-Personalized		Personalized		Non-Personalized		Personalized	
		FlanT5-XXL	GPT-3.5	FlanT5-XXL	GPT-3.5	FlanT5-XXL	GPT-3.5	FlanT5-XXL	GPT-3.5
LaMP-1: Personalized Citation Identification	Accuracy	0.520	0.541	0.699	0.695	0.502	0.508	0.636	0.634
LaMP-2: Personalized News Categorization	Accuracy	0.581	0.594	0.617	0.643	0.632	0.691	0.600	0.684
	F1	0.475	0.488	0.512	0.563	0.520	0.536	0.519	0.551
LaMP-3: Personalized Product Rating	MAE	0.344	0.706	0.267	0.620	0.333	0.677	0.299	0.603
	RMSE	0.650	0.972	0.552	1.049	0.650	0.948	0.616	1.002
LaMP-4: Personalized News Headline Generation	ROUGE-1	0.163	0.136	0.182	0.150	0.176	0.146	0.188	0.158
	ROUGE-L	0.147	0.119	0.167	0.133	0.160	0.128	0.172	0.140
LaMP-5: Personalized Scholarly Title Generation	ROUGE-1	0.442	0.387	0.450	0.390	0.471	0.424	0.483	0.425
	ROUGE-L	0.400	0.329	0.411	0.329	0.422	0.355	0.433	0.351
LaMP-6: Personalized Email Subject Generation	ROUGE-1	0.362	-	0.482	-	0.335	-	0.401	-
	ROUGE-L	0.343	-	0.471	-	0.319	-	0.387	-
LaMP-7: Personalized Tweet Paraphrasing	ROUGE-1	0.453	0.399	0.448	0.390	0.448	0.390	0.440	0.382
	ROUGE-L	0.395	0.336	0.394	0.322	0.396	0.330	0.389	0.318

Table 3: The zero-shot personalized text classification and generation results on the test set of user- and time-based separation settings. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3U and LaMP-3T datasets. For personalized models, the tuned retriever based on the validation performance was selected. The results on the validation set are reported in Tables 6 and 8 in Appendix E.

The results suggest that increasing the number of retrieved items leads to improved performance in downstream tasks. However, certain tasks experience a decline in performance under these conditions. Given the finite context size of language models, exploring approaches to generate a unified prompt from multiple user entries appears to be a promising avenue for future investigation.

The Impact of Tuning Retriever Parameters on the End-to-End Performance. Based on the performance on the validation set for each dataset, we tuned two parameters for each dataset: (1) the retrieval model (BM25 vs. Contriever vs. Recency), and (2) the number of retrieved items (k). For parameter tuning, we used the following metrics: Accuracy for LaMP-1 and LaMP-2, MAE for LaMP-3, and ROUGE-1 for all text generation tasks. The results for this tuned model are presented in the last column of Table 1 and Table 2. As expected, the tuned model outperforms the other models on all datasets.

4.3 Zero-Shot Personalized Classification and Generation Results

In view of the widespread adoption of employing large-scale language models with no fine-tuning in contemporary research, we conduct an evaluation of two such models on our benchmark dataset. In particular, we leverage GPT 3.5 (alias gpt-3.5-turbo or ChatGPT⁵ and FlanT5-XXL

(Chung et al., 2022). FlanT5-XXL comprises 11B parameters, however, the size of GPT-3.5 is unknown (GPT3 consists of 175B parameters). For evaluation, we provide each model with the inputs corresponding to individual tasks and assess their performance based on the generated outputs. In the context of classification tasks, if the produced output does not correspond to a valid class, we resort to calculating the similarity between each class label and the generated output utilizing BERTScore (Zhang* et al., 2020). Subsequently, we assign the most similar label to the generated output as the corresponding output for the given input. To shed light on this, GPT-3.5 generated out-of-the-label predictions 8%, 4%, 6%, 4%, 2%, and 4% of the time for the LaMP-1U, LaMP-1T, LaMP-2U, LaMP-2T, LaMP-3U, and LaMP-3T tasks, respectively. On the other hand, FlanT5-XXL predictions are consistently among the questioned labels.

Table 3 shows the result of LLMs on this benchmark in a zero-shot scenario. The results show that, except for the Personalized Tweet Paraphrasing task, using the user’s profile with LLMs improves their performance on this benchmark in a zero-shot setting. The outcomes in Tables 1 and 2 show the results for FlanT5-base, a 250M parameter model, fine-tuned on each task. Table 3 presents the zero-shot application of LLMs. These findings indicate that fine-tuning smaller models on downstream tasks leads to enhanced performance in comparison to zero-shot performance of LLMs.

Finally, it is crucial to highlight that the observed

⁵<https://openai.com/blog/chatgpt>

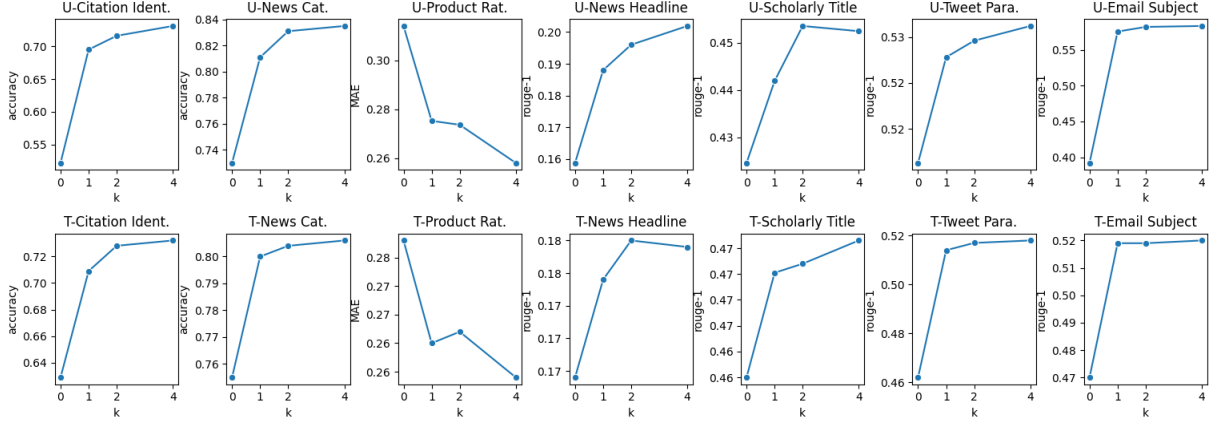


Figure 2: The performance on downstream tasks using the best retriever for each task from Tables 1 and 2 with different numbers of retrieved entries from the user profile. It can be seen that increasing the number of profile entries usually increases the performance of the downstream tasks.

outcomes, which indicate superior performance of FlanT5-XXL over GPT-3.5, should not be construed as an inherent deficiency of the latter model. The efficacy of large language models (LLMs) is extensively contingent upon the caliber and configuration of the input prompts employed. It is worth noting that prompt engineering, which plays a significant role in determining the performance of LLMs in the given tasks, is not the central objective of this study. Consequently, any disparities in performance must be evaluated in light of this contextual information.

5 Research Problems Enabled by LaMP

In this section, we explain several open problems that the LaMP benchmark can play an essential role in moving through cracking them. The open problems include and are not limited to:

Research on Prompting Language Models for Personalization. One straightforward approach to integrating user profiles into language models is through the use of hard prompts (Brown et al., 2020). However, the limited context size of LMs makes it impractical and/or expensive to augment them with lengthy user profile entries. Studying various prompts for personalization purposes would be interesting. Generating personalization prompts based on the user profile, instead of the employing the retrieved entries would be an alternative solution for future exploration. Additionally, employing soft prompts (Lester et al., 2021) can be further explored in the context of personalizing language model outputs by learning personalized soft prompt vectors appears to be a more effective

solution to this issue.

Evaluation of Personalized Text Generation.

The commonly used evaluation metrics for text generation, whether syntactical (Lin, 2004; Banerjee and Lavie, 2005; Papineni et al., 2002) or semantical (Zhang* et al., 2020), do not incorporate the user into their evaluation process. Consequently, such metrics may not be entirely suitable for evaluating personalized text generation problems. Therefore, exploring and developing new evaluation metrics that better account for the user’s preferences and characteristics can significantly benefit this area of research.

Learning to Retrieve from User Profiles.

Learning to rank has been widely explored in various retrieval scenarios. Optimizing ranking models that select personalized entries for the sake of personalized text classification and/or generation would be a potentially impactful research direction.

6 Related Work

Personalization has been well studied for information access problems, with the organization of the Netflix Challenge and its associated datasets representing an important driver of academic focus on personalization (Konstan and Terveen, 2021). It also represents an important element of large-scale industry recommender systems (Davidson et al., 2010; Das et al., 2007; Xu et al., 2022) and has also been extensively studied for search applications (Bennett et al., 2012; Dumais, 2016; Croft et al., 2001; Tabrizi et al., 2018; Zeng et al., 2023), in contexts ranging from query auto-completion

(Jaech and Ostendorf, 2018) to collaborative personalized search (Xue et al., 2009). We refer readers to Rafeian and Yoganarasimhan (2023) for an overview of this line of work. Here, we cover work on personalization in NLP, focusing on the datasets available for research.

Personalization has been examined extensively for dialogue agents (Wu et al., 2021; Zhang et al., 2018; Mazaré et al., 2018). Compared to other NLP tasks, this focus likely stems from the importance of tailoring dialogue to users and conditioning generated utterances on specific personas. Given the lack of real conversational data, this work has constructed dialogue data for users by promoting crowd-workers to author dialogues based on specific personas (Zhang et al., 2018), and through extracting user attributes and utterances from Reddit (Mazaré et al., 2018; Wu et al., 2021) and Weibo (Zhong et al., 2022; Qian et al., 2021). To leverage more realistic conversational data, recent work of Vincent et al. (2023) annotate a dataset of movie dialogues with narrative character personas and posit the potential for using large-language models for dialogue generation conditioned on these personas. Besides exploring text generation for dialogues, other work has also leveraged publicly available reviews and recipes to explore personalization for review (Li and Tuzhilin, 2019) and recipe generation tasks (Majumder et al., 2019). And Wuebker et al. (2018) explore parameter efficient models for personalized translation models with a non-public dataset. Finally, Ao et al. (2021) presents a personalized headline generation dataset constructed from realistic user interaction data on Microsoft News. This is closely related to the headline generation task of LaMP, which focuses on personalization for *authors* rather than readers. More broadly, LaMP presents resources for the tasks which have seen lesser attention than those based on dialogue – expanding the underexplored space of personalizing text generation systems (Dudy et al., 2021).

While a body of work has focused on user-facing applications, others have explored personalization for more fundamental problems in language modeling. This body of work has leveraged openly available user data on Reddit (Welch et al., 2022), Facebook, Twitter (Soni et al., 2022), and other blogging sites (King and Cook, 2020). Besides pre-training language models for personalization, Soni et al. (2022) also explore applying a personalized language model for downstream tasks in stance

classification and demographic inference. Similarly, other work has explored personalized sentiment prediction tasks on publicly available Yelp and IMDB data (Mireshghallah et al., 2022; Zhong et al., 2021) – this work bears a resemblance to the personalized product rating task in LaMP and ties back to rating prediction tasks explored in recommendation tasks. Finally, Plepi et al. (2022) examines the application of personalization methods to modeling annotators in a classification task reliant on modeling social norms – making an important connection between personalization and an emerging body of work on accommodating human label variation in NLP research (Rottger et al., 2022; Gordon et al., 2022; Plank, 2022).

7 Conclusion

This paper presents a novel benchmark named LaMP for training and evaluating language models for personalized text classification and generation tasks. LaMP consists of seven datasets: three classification (including two categorical and one ordinal) and four generation datasets. To establish baseline performance, we perform extensive experiments using various language models and retrieval techniques for selecting user profile entries for producing personalized prompts. Lastly, we report the performance of prominent large-scale language models on this benchmark.

Limitations

A shortcoming of the present study is that, although all tasks are designed to assess language models’ proficiency in personalization, certain tasks could be better grounded in realistic scenarios and real-world applications. For instance, framing the Personalized Citation Identification task as a binary classification problem might not accurately represent real-world situations, where individuals generally need to interact with a more extensive array of articles. Additionally, while Personalized Product Rating is intrinsically linked to predicting user satisfaction, the approach may not be entirely realistic, as reviews in real-world contexts are often accompanied by a numerical rating, rendering direct score prediction less relevant. That being said, this benchmark creates an environment for evaluating the abilities of models in producing personalized outputs.

Furthermore, the data used for creating the LaMP benchmark mostly consists of publicly avail-

able data on the Web, e.g., public tweets, scholarly articles, news articles, and product reviews. We should take this into consideration that some of this data may have been observed by the large language models during their pretraining process. Therefore, they may even perform poorer in unseen cases compared to what we observe from the results on most the LaMP datasets. For this reason, we included the Avocado dataset for Personalized Email Subject Generation as this is not publicly available on the Web and we expect that language models do not use this dataset for pretraining given the restrictions on the data usage agreement.

Moreover, this benchmark operates under the assumption that there are no common users among the train, validation, and test sets. In contrast, real-world scenarios often employ a warm-start approach, where the evaluation of a user using a model trained with their prior history rather than the history of other users. This method could potentially lead to more accurate personalization and model performance, thus highlighting another area where the benchmark may not fully capture real-world challenges.

Acknowledgments

This work was supported in part by the Center for Intelligent Information Retrieval, in part by NSF grant #2143434, in part by the Office of Naval Research contract number N000142212688, and in part by Lowe’s. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect those of the sponsors.

References

- Xiang Ao, Xiting Wang, Ling Luo, Ying Qiao, Qing He, and Xing Xie. 2021. [PENS: A dataset and generic framework for personalized news headline generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 82–92, Online. Association for Computational Linguistics.
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Paul N. Bennett, Ryen W. White, Wei Chu, Susan T. Dumais, Peter Bailey, Fedor Borisjuk, and Xiaoyuan Cui. 2012. [Modeling the impact of short- and long-term behavior on search personalization](#). In *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’12*, page 185–194, New York, NY, USA. Association for Computing Machinery.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#).
- Bruce W. Croft, Stephen Cronen-Townsend, and Victor Lavrenko. 2001. [Relevance feedback and personalization: A language modeling perspective](#). In *DELOS Workshop: Personalisation and Recommender Systems in Digital Libraries*.
- Abhinandan S. Das, Mayur Datar, Ashutosh Garg, and Shyam Rajaram. 2007. [Google news personalization: Scalable online collaborative filtering](#). In *Proceedings of the 16th International Conference on World Wide Web, WWW ’07*, page 271–280, New York, NY, USA. Association for Computing Machinery.
- James Davidson, Benjamin Liebald, Junning Liu, Palash Nandy, Taylor Van Vleet, Ullas Gargi, Sujoy Gupta, Yu He, Mike Lambert, Blake Livingston, and Dasarathi Sampath. 2010. [The youtube video recommendation system](#). In *Proceedings of the Fourth ACM Conference on Recommender Systems, RecSys ’10*, page 293–296, New York, NY, USA. Association for Computing Machinery.
- Shiran Dudy, Steven Bedrick, and Bonnie Webber. 2021. [Refocusing on relevance: Personalization in NLG](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*,

- pages 5190–5202, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Susan T. Dumais. 2016. [Personalized search: Potential and pitfalls](#). In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, CIKM '16, page 689, New York, NY, USA. Association for Computing Machinery.
- Peter S. Fader, Bruce G.S. Hardie, and Ka Lok Lee. 2005. [Rfm and clv: Using iso-value curves for customer base analysis](#). *Journal of Marketing Research*, 42(4):415–430.
- Lucie Flek. 2020. [Returning the N to NLP: Towards contextually personalized classification models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7828–7838, Online. Association for Computational Linguistics.
- Andrew Fowler, Kurt Partridge, Ciprian Chelba, Xiaojun Bi, Tom Ouyang, and Shumin Zhai. 2015. [Effects of language modeling and its personalization on touchscreen typing performance](#). In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, CHI '15, page 649–658, New York, NY, USA. Association for Computing Machinery.
- Markus Freitag and Yaser Al-Onaizan. 2017. [Beam search strategies for neural machine translation](#). In *Proceedings of the First Workshop on Neural Machine Translation*, pages 56–60, Vancouver. Association for Computational Linguistics.
- Sebastian Gehrmann, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu, Dipanjan Das, Kaustubh Dhole, Wanyu Du, Esin Durmus, Ondřej Dušek, Chris Chenyene Emezue, Varun Gangal, Cristina Garbacea, Tatsunori Hashimoto, Yufang Hou, Yacine Jernite, Harsh Jhamtani, Yangfeng Ji, Shailza Jolly, Mihir Kale, Dhruv Kumar, Faisal Ladhak, Aman Madaan, Mounica Maddela, Khyati Mahajan, Saad Mahamood, Bodhisattwa Prasad Majumder, Pedro Henrique Martins, Angelina McMillan-Major, Simon Mille, Emiel van Miltenburg, Moin Nadeem, Shashi Narayan, Vitaly Nikolaev, Andre Niyongabo Rubungo, Salomey Osei, Ankur Parikh, Laura Perez-Beltrachini, Niranjana Ramesh Rao, Vikas Raunak, Juan Diego Rodriguez, Sashank Santhanam, João Sedoc, Thibault Sellam, Samira Shaikh, Anastasia Shimorina, Marco Antonio Sobrevilla Cabeza, Hendrik Strobelt, Nishant Subramani, Wei Xu, Diyi Yang, Akhila Yerukola, and Jiawei Zhou. 2021. [The GEM benchmark: Natural language generation, its evaluation and metrics](#). In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 96–120, Online. Association for Computational Linguistics.
- Alec Go, Richa Bhayani, and Lei Huang. 2009. [Twitter sentiment classification using distant supervision](#).
- Mitchell L. Gordon, Michelle S. Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S. Bernstein. 2022. [Jury learning: Integrating dissenting voices into machine learning models](#). In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22, New York, NY, USA. Association for Computing Machinery.
- Xiaolei Huang, Lucie Flek, Franck Dernoncourt, Charles Welch, Silvio Amir, Ramit Sawhney, and Diyi Yang. 2022. [UserNLP'22: 2022 international workshop on user-centered natural language processing](#). In *Companion Proceedings of the Web Conference 2022*, WWW '22, page 1176–1177, New York, NY, USA. Association for Computing Machinery.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. [Unsupervised dense information retrieval with contrastive learning](#). *Transactions on Machine Learning Research*.
- Aaron Jaech and Mari Ostendorf. 2018. [Personalized language model for query auto-completion](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 700–705, Melbourne, Australia. Association for Computational Linguistics.
- Milton King and Paul Cook. 2020. [Evaluating approaches to personalizing language models](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2461–2469, Marseille, France. European Language Resources Association.
- Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A Hale. 2023. Personalisation within bounds: A risk taxonomy and policy framework for the alignment of large language models with personalised feedback. *arXiv preprint arXiv:2303.05453*.
- Joseph Konstan and Loren Terveen. 2021. [Human-centered recommender systems: Origins, advances, challenges, and opportunities](#). *AI Magazine*, 42(3):31–42.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Pan Li and Alexander Tuzhilin. 2019. [Towards controllable and personalized review generation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3237–3245, Hong Kong, China. Association for Computational Linguistics.

- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. 2019. [Generating personalized recipes from historical user preferences](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5976–5982, Hong Kong, China. Association for Computational Linguistics.
- Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. [Training millions of personalized dialogue agents](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2775–2779, Brussels, Belgium. Association for Computational Linguistics.
- Fatemehsadat Mireshghallah, Vaishnavi Shrivastava, Milad Shokouhi, Taylor Berg-Kirkpatrick, Robert Sim, and Dimitrios Dimitriadis. 2022. [UserIdentifier: Implicit user representations for simple and effective personalized sentiment analysis](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3449–3456, Seattle, United States. Association for Computational Linguistics.
- Rishabh Misra. 2022. News category dataset. *arXiv preprint arXiv:2209.11429*.
- Rishabh Misra and Jigyasa Grover. 2021. *Sculpting Data for ML: The first act of Machine Learning*.
- Maxim Naumov, Dheevatsa Mudigere, Hao-Jun Michael Shi, Jianyu Huang, Narayanan Sundaraman, Jongsoo Park, Xiaodong Wang, Udit Gupta, Carole-Jean Wu, Alisson G. Azzolini, Dmytro Dzhulgakov, Andrey Malleevich, Ilia Cherniavskii, Yinghai Lu, Raghuraman Krishnamoorthi, Ansha Yu, Volodymyr Kondratenko, Stephanie Pereira, Xianjie Chen, Wenlin Chen, Vijay Rao, Bill Jia, Liang Xiong, and Misha Smelyanskiy. 2019. [Deep learning recommendation model for personalization and recommendation systems](#).
- Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. [Justifying recommendations using distantly-labeled reviews and fine-grained aspects](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 188–197, Hong Kong, China. Association for Computational Linguistics.
- Oard, Douglas, Webber, William, Kirsch, David A., and Golitsynskiy, Sergey. 2015. [Avocado research email collection](#).
- OpenAI. 2023. [Gpt-4 technical report](#).
- Sheena Panthaplackel, Adrian Benton, and Mark Dredze. 2022. [Updated headline generation: Creating updated summaries for evolving news stories](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6438–6461, Dublin, Ireland. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. [KILT: a benchmark for knowledge intensive language tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544, Online. Association for Computational Linguistics.
- Barbara Plank. 2022. [The “problem” of human label variation: On ground truth in data, modeling and evaluation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10671–10682, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Joan Plepi, Béla Neuendorf, Lucie Flek, and Charles Welch. 2022. [Unifying data perspectivism and personalization: An application to social norms](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7391–7402, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Hongjin Qian, Xiaohe Li, Hanxun Zhong, Yu Guo, Yueyuan Ma, Yutao Zhu, Zhanliang Liu, Zhicheng Dou, and Ji-Rong Wen. 2021. [Pchatbot: A large-scale dataset for personalized chatbot](#). In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’21*, page 2470–2477, New York, NY, USA. Association for Computing Machinery.
- Omid Rafieian and Hema Yoganarasimhan. 2023. Ai and personalization. *Artificial Intelligence in Marketing*, pages 77–102.
- Werner J. Reinartz and V. Kumar. 2000. [On the profitability of long-life customers in a noncontractual setting: An empirical investigation and implications for marketing](#). *Journal of Marketing*, 64(4):17–35.

- Werner J. Reinartz and V. Kumar. 2003. [The impact of customer relationship characteristics on profitable lifetime duration](#). *Journal of Marketing*, 67(1):77–99.
- Stephen Robertson, S. Walker, S. Jones, M. M. Hancock-Beaulieu, and M. Gatford. 1995. Okapi at trec-3. In *Proceedings of the Third Text REtrieval Conference, TREC-3*, pages 109–126. Gaithersburg, MD: NIST.
- Paul Rottger, Bertie Vidgen, Dirk Hovy, and Janet Pierrehumbert. 2022. [Two contrasting data annotation paradigms for subjective NLP tasks](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 175–190, Seattle, United States. Association for Computational Linguistics.
- Nikita Soni, Matthew Matero, Niranjana Balasubramanian, and H. Andrew Schwartz. 2022. [Human language modeling](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 622–636, Dublin, Ireland. Association for Computational Linguistics.
- Shayan A. Tabrizi, Azadeh Shakery, Hamed Zamani, and Mohammad Ali Tavallaei. 2018. [Person: Personalized information retrieval evaluation based on citation networks](#). *Information Processing & Management*, 54(4):630–656.
- Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. 2008. Arnetminer: Extraction and mining of academic social networks. In *KDD’08*, pages 990–998.
- Sebastian Vincent, Rowanne Sumner, Alice Dowek, Charlotte Blundell, Emily Preston, Chris Bayliss, Chris Oakley, and Carolina Scarton. 2023. Personalised language modelling of screen characters using rich metadata annotations. *arXiv preprint arXiv:2303.16618*.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. [Superglue: A stickier benchmark for general-purpose language understanding systems](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Charles Welch, Chenxi Gu, Jonathan K. Kummerfeld, Veronica Perez-Rosas, and Rada Mihalcea. 2022. [Leveraging similar users for personalized language modeling with limited data](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1742–1752, Dublin, Ireland. Association for Computational Linguistics.
- Yuwei Wu, Xuezhe Ma, and Diyi Yang. 2021. [Personalized response generation via generative split memory network](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1956–1970, Online. Association for Computational Linguistics.
- Joern Wuebker, Patrick Simianer, and John DeNero. 2018. [Compact personalized models for neural machine translation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 881–886, Brussels, Belgium. Association for Computational Linguistics.
- Jiajing Xu, Andrew Zhai, and Charles Rosenberg. 2022. [Rethinking personalized ranking at pinterest: An end-to-end approach](#). In *Proceedings of the 16th ACM Conference on Recommender Systems, RecSys ’22*, page 502–505, New York, NY, USA. Association for Computing Machinery.
- Gui-Rong Xue, Jie Han, Yong Yu, and Qiang Yang. 2009. [User language model for collaborative personalized search](#). *ACM Trans. Inf. Syst.*, 27(2).
- Kosuke Yamada, Yuta Hitomi, Hideaki Tamori, Ryoei Sasano, Naoaki Okazaki, Kentaro Inui, and Koichi Takeda. 2021. [Transformer-based lexically constrained headline generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4085–4090, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hansi Zeng, Surya Kallumadi, Zaid Alibadi, Rodrigo Nogueira, and Hamed Zamani. 2023. A personalized dense retrieval framework for unified information access. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’23*, New York, NY, USA. Association for Computing Machinery.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.
- Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Hanxun Zhong, Zhicheng Dou, Yutao Zhu, Hongjin Qian, and Ji-Rong Wen. 2022. [Less is more: Learning to refine dialogue history for personalized dialogue generation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5808–5820, Seattle, United States. Association for Computational Linguistics.

Wanjun Zhong, Duyu Tang, Jiahai Wang, Jian Yin, and Nan Duan. 2021. [UserAdapter: Few-shot user learning in sentiment analysis](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1484–1488, Online. Association for Computational Linguistics.

Jianing Zhou and Suma Bhat. 2021. [Paraphrase generation: A survey of the state of the art](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5075–5086, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

A Data Creation Details for Tasks in the LaMP benchmark

This section explains the details behind creating inputs, outputs, and profile entries for each task in the LaMP benchmark.

A.1 User-based Separation Setting

Personalized Citation Identification. To generate data samples, we leverage the Citation Network Dataset (V14) (Tang et al., 2008), which comprises information on scientific papers, authors, and citations. We select all papers from this dataset that meet the following criteria: 1) they are written in English, 2) they contain at least one reference and one author, and 3) they include an abstract. Subsequently, we group papers based on their authors and only consider authors who have written at least 50 papers. For each author, we randomly select one of their papers and one of its cited references. For negative document selection, we randomly choose one of the first author’s co-authors and one of the papers they have cited in one of their papers, which has not been cited by the first author. If no such author exists, we randomly select an author and repeat this process. Finally, we construct the input, output, and profile of the generated samples for this task, employing the template depicted in Figure 3 in Appendix C. After creating the samples for all users, we divide users into the train, validation, and test sets for this task.

Personalized News Categorization. To construct our dataset for this task, we leverage the news categorization dataset (Misra, 2022; Misra and Grover, 2021) obtained from the HuffPost website⁶. However, we filter out some categories and merge similar ones to form a more concise set of categories, as described in Appendix B. Next, we group articles by their author, taking only the first author in cases where there are multiple authors, and retain only authors with a minimum of four articles who have published in at least three different categories. Then, we partition the authors into training, validation, and test sets. For each article of an author, we use the article as input, the article’s category as output, and the remaining articles of the same author as the user profile for that sample, following the template shown in Figure 3 in Appendix C. Finally, we randomly select 50% of the generated samples for each user in training, validation, and test sets, and add them to the samples of the corresponding set.

Personalized Product Rating. In this task, we create our dataset by leveraging the Amazon Reviews Dataset (Ni et al., 2019). We filtered out users (i.e., amazon customers who have written reviews) who have written less than 100 and the 1% users with the most reviews as outliers. Since the Amazon Reviews dataset is quite extensive, we randomly sampled a subset of users from the dataset, which we then split into training, validation, and testing sets.

To construct the input-output pairs for our task, for each user, we randomly select one of their reviews as the input to the task and use their other reviews as their profile. Specifically, we use the profile to capture the author’s writing style, preferences, and tendencies. In this setup, the user’s score for the input review serves as the ground truth output for our task. To gain a better understanding of the input, output, and profile, refer to Figure 3 in Appendix C.

Personalized News Headline Generation. To construct our dataset for this task, we leverage the News Categorization dataset (Misra, 2022; Misra and Grover, 2021) from the HuffPost website⁷. The dataset provides author information for each article and is used to group articles by their respective authors. Similar to Section 3.2, we filtered out the

⁶<https://www.huffpost.com/>

⁷<https://www.huffpost.com/>

authors with less than four articles. In cases where an article has multiple authors, we assign it only to the first author.

We then randomly split the authors into training, validation, and test sets. For each author in each set, we create input-output pairs by selecting each article as the input, the headline of the article as the output, and the remaining articles written by the same author as their profile. This setup aims to capture the author’s writing style, preferences, and tendencies, which can be leveraged to generate headlines that align with their interests. An example of this setup is presented in Figure 3 in Appendix C. Finally, to ensure a diverse and representative dataset, we randomly select 50% of the created samples for each author and add them to the user’s corresponding set.

Personalized Scholarly Title Generation. Similar to Section 3.1, we leverage the Citation Network Dataset (V14) (Tang et al., 2008) that includes information about scientific papers, authors, and citations to construct our dataset. We only kept the papers that meet the following criteria: 1) written in English, 2) have at least one reference and one author, and 3) have an abstract. Then, we group papers by their authors and only consider authors who have published at least 50 papers. For each author, we randomly choose one of their papers and use its abstract as input, its title as output, and the remaining papers as the author’s profile. Figure 3 in Appendix C illustrates the input format for this task. After creating the samples for all users, we divide users into the train, validation, and test sets for this task.

Personalized Email Subject Generation. In this study, we adopt the Avocado Research Email Collection (Oard, Douglas et al., 2015) as the primary dataset for our task. To curate the dataset, we first perform a filtering step where we exclude emails with subject lengths of fewer than five words and content lengths of fewer than 30 words. Next, we group the emails based on their sender’s email address, retaining only those from users with email frequencies ranging between 10 to 200 emails. We further divide the users into distinct training, validation, and test sets to ensure that our model generalizes well to unseen data. To generate training examples for each user, we create input-output pairs by considering each email as the input and the corresponding email subject as the output. We sup-

plement these pairs with other emails written by the same user as their profile, as shown in Figure 3 in Appendix C. We ensure that our dataset is diverse and representative by randomly selecting 50% of the curated samples for each user and adding them to their respective sets.

Personalized Tweet Paraphrasing. In this task, we utilize the Sentiment140 dataset (Go et al., 2009) as our tweet collection set. To ensure that the collected tweets are of adequate length, we only retain tweets containing at least 10 words. We then group the tweets based on the user ID and filter out users with fewer than 10 tweets. Subsequently, we randomly select one tweet from each user profile and use it as input to ChatGPT (i.e., gpt3.5-turbo)⁸ to generate a paraphrased version. The generated paraphrase is then utilized as the input to our NLP task, with the original tweet serving as the corresponding output. The remaining tweets of the user constitute the user’s profile, excluding the one selected as input. Figure 3 in Appendix C provides an overview of the input-output-profile template for our proposed task. After creating the samples for all users, we divide users into the train, validation, and test sets for this task.

A.2 Time-based Separation Setting

Personalized Citation Identification. To generate data samples, we leverage the Citation Network Dataset (V14) (Tang et al., 2008), which comprises information on scientific papers, authors, and citations. We select all papers from this dataset that meet the following criteria: 1) they are written in English, 2) they contain at least one reference and one author, and 3) they include an abstract. Subsequently, we group papers based on their authors and only consider authors who have written at least 50 papers. We divide each author’s papers based on the publication year into three groups chronologically: 1) profile papers, 2) train papers, 3) validation papers, and 4) test papers, where the order of groups shows the flow of time. Each train, validation, and test paper set in this task consists of only one paper. Then, for each paper in the train, validation, and test sets for the user, we select each paper and one of its cited references. For negative document selection, we randomly choose one of the first author’s co-authors and one of the papers they have cited in one of their papers, which has not been cited by the first author. If no such author

⁸<https://openai.com/blog/chatgpt>

exists, we randomly select an author and repeat this process. Finally, we construct the input, output, and profile of the generated samples for this task, employing the template depicted in Figure 3 in Appendix C. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of validation and test sets to create the final sets for the task.

Personalized News Categorization. To construct our dataset for this task, we leverage the news categorization dataset (Misra, 2022; Misra and Grover, 2021) obtained from the HuffPost website⁹. However, we filter out some categories and merge similar ones to form a more concise set of categories, as described in Appendix B. Next, we group articles by their author, taking only the first author in cases where there are multiple authors, and retain only authors with a minimum of 10 articles who have published in at least three different categories. We divide each author’s articles based on the publishing date into three groups chronologically: 1) profile articles, 2) train articles, 3) validation articles, and 4) test articles, where the order of groups shows the flow of time. Each train, validation, and test articles set in this task consists of 20%, 10%, and 10% articles, respectively. Then, for each article in the train, validation, and test sets for the user, we use the article as input, the article’s category as output, and the profile articles as the profile for that sample, following the template shown in Figure 3 in Appendix C. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of validation and test sets to create the final sets for the task.

Personalized Product Rating. In this task, we create our dataset by leveraging the Amazon Reviews Dataset (Ni et al., 2019). We filtered out users (i.e., amazon customers who have written reviews) who have written less than 100 and the 1% users with the most reviews as outliers. Since the Amazon Reviews dataset is quite extensive, we randomly sampled a subset of users from the dataset. We divide each user’s reviews based on the review date into three groups chronologically: 1) profile reviews, 2) train reviews, 3) validation reviews, and 4) test reviews, where the order of groups shows

the flow of time. Each train, validation, and test reviews set in this task consists of only one review.

To construct the input-output pairs for our task, for each user, we select their reviews in each of train, validation, and test sets as the input to the task and use the profile reviews as their profile. Specifically, we use the profile to capture the author’s writing style, preferences, and tendencies. Additionally, the user’s score for the input review serves as the ground truth output for our task. To gain a better understanding of the input, output, and profile, refer to Figure 3 in Appendix C. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of validation and test sets to create the final sets for the task.

Personalized News Headline Generation. To construct our dataset for this task, we leverage the News Categorization dataset (Misra, 2022; Misra and Grover, 2021) from the HuffPost website¹⁰. The dataset provides author information for each article and is used to group articles by their respective authors. We filtered out the authors with less than ten articles. In cases where an article has multiple authors, we assign it only to the first author. We divide each author’s articles based on the publishing date into three groups chronologically: 1) profile articles, 2) train articles, 3) validation articles, and 4) test articles, where the order of groups shows the flow of time. Each train, validation, and test articles set in this task consists of 20%, 10%, and 10% articles, respectively. Then, for each article in the train, validation, and test sets for the user, we create input-output pairs by selecting each article as the input, the headline of the article as the output, and the profile articles written by the same author as their profile. This setup aims to capture the author’s writing style, preferences, and tendencies, which can be leveraged to generate headlines that align with their interests. An example of this setup is presented in Figure 3 in Appendix C. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of validation and test sets to create the final sets for the task.

Personalized Scholarly Title Generation. Similar to Section 3.1, we leverage the Citation Net-

⁹<https://www.huffpost.com/>

¹⁰<https://www.huffpost.com/>

work Dataset (V14) (Tang et al., 2008) that includes information about scientific papers, authors, and citations to construct our dataset. We only kept the papers that meet the following criteria: 1) written in English, 2) have at least one reference and one author, and 3) have an abstract. Then, we group papers by their authors and only consider authors who have published at least 50 papers. We divide each author’s papers based on the publication year into three groups chronologically: 1) profile papers, 2) train papers, 3) validation papers, and 4) test papers, where the order of groups shows the flow of time. Each train, validation, and test paper set in this task consists of only one paper. Then, for each paper in the train, validation, and test sets for the user, we use its abstract as input, its title as output, and the profile papers as the author’s profile. Figure 3 in Appendix C illustrates the input format for this task. After creating the samples for all users, we divide users into the train, validation, and test sets for this task. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of validation and test sets to create the final sets for the task.

Personalized Email Subject Generation. In this study, we adopt the Avocado Research Email Collection (Oard, Douglas et al., 2015) as the primary dataset for our task. To curate the dataset, we first perform a filtering step where we exclude emails with subject lengths of fewer than five words and content lengths of fewer than 30 words. Next, we group the emails based on their sender’s email address, retaining only those from users with email frequencies ranging between 10 to 200 emails. We divide each user’s emails based on the publishing date into three groups chronologically: 1) profile emails, 2) train emails, 3) validation emails, and 4) test emails, where the order of groups shows the flow of time. Each train, validation, and test emails set in this task consists of 20%, 10%, and 10% articles, respectively. Then, for each article in the train, validation, and test sets for the user, we create input-output pairs by considering each email as the input and the corresponding email subject as the output. We supplement these pairs with profile emails written by the same user as their profile, as shown in Figure 3 in Appendix C. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of

validation and test sets to create the final sets for the task.

Personalized Tweet Paraphrasing. In this task, we utilize the Sentiment140 dataset (Go et al., 2009) as our tweet collection set. To ensure that the collected tweets are of adequate length, we only retain tweets containing at least 10 words. We then group the tweets based on the user ID and filter out users with fewer than 10 tweets. We divide each user’s tweets based on the publication year into three groups chronologically: 1) profile tweets, 2) train tweets, 3) validation tweets, and 4) test tweets, where the order of groups shows the flow of time. Each train, validation, and test tweet set in this task consists of only one paper. Then, for each tweet in the train, validation, and test sets for the user, we use it as input to ChatGPT (i.e., gpt3.5-turbo)¹¹ to generate a paraphrased version. The generated paraphrase is then utilized as the input to our NLP task, with the original tweet serving as the corresponding output. The profile tweets of the user constitute the user’s profile. Figure 3 in Appendix C provides an overview of the input-output-profile template for our proposed task. It should be noted that after creating all the samples in the train, validation, and test sets for all the users and aggregating them, we randomly select a subset of validation and test sets to create the final sets for the task.

B News Categorization Dataset Category Merging Strategy

The original news categorization dataset contains the following categories: POLITICS, ENTERTAINMENT, HEALTHY LIVING, BUSINESS, SPORTS, COMEDY, PARENTS, WOMEN, CRIME, RELIGION, STYLE, TRAVEL, ARTS, SCIENCE, ARTS & CULTURE, TECH, COLLEGE, EDUCATION, STYLE & BEAUTY, CULTURE & ARTS, FOOD & DRINK. We decided only to keep the following categories based on their distribution: POLITICS, ENTERTAINMENT, HEALTHY LIVING, BUSINESS, SPORTS, COMEDY, PARENTS, WOMEN, CRIME, RELIGION, STYLE, TRAVEL, ARTS, SCIENCE, ARTS & CULTURE, TECH, COLLEGE, EDUCATION, STYLE & BEAUTY, CULTURE & ARTS, FOOD & DRINK.

Next, based on the similarity of different cate-

¹¹<https://openai.com/blog/chatgpt>

gories, we merge the following categories into the same category:

- “ENTERTAINMENT” and “COMEDY” are merged as “ENTERTAINMENT”.
- “STYLE” and “STYLE & BEAUTY” are merged as “STYLE & BEAUTY”.
- “ARTS”, “CULTURE & ARTS”, and “ARTS & CULTURE” are merged as “CULTURE & ARTS”.
- “COLLEGE” and “EDUCATION” are merged as “EDUCATION”.
- “SCIENCE” and “TECH” are merged as “SCIENCE & TECHNOLOGY”.

C Samples of the Tasks Introduced in the LaMP Benchmark

As mentioned earlier, LaMP proposes seven tasks to evaluate language model personalization. In order to create the data points, we use just a carefully designed template for each task. Figure 3 depicts a sample and template for each task in LaMP. Generally, each sample in each task has an input and output accompanied by a profile consisting of several entries about the user, helping the model to produce personalized results for the user. While the profile entries in the same task have a similar structure, the structure varies between tasks. For example, Figure 3 shows that the profile for Personalized Product Rating comprised of documents with text and score sections, while the profile entries in Personalized Scholarly Title Generation have abstract and title attributes.

D Prompts Used for Adding User Profile to the Language Model’s Input

In order to use multiple entries from the user profile to personalize the language model’s input, we construct task-specific prompts using the templates and instructions in Table 4.

The prompt creation consists of two stages: 1) Per Profile Entry Prompt (PPEP) creation and 2) Aggregated Input Prompt (AIP) creation. In the first stage, following the instructions in Table 4, we create a prompt for each profile entry. In the second stage, following the instructions in Table 4, we combine the PPEP prompts into a single prompt to be fed to the language model. It should

be noted that due to the limited context size of language models, we need to trim the PPEP prompts. More accurately, considering k prompts need to be merged and that the maximum capacity for the task input is \bar{L} and the maximum context size of the language model is L , we let each PPEP occupy $\frac{L-\bar{L}}{k}$ tokens in the language model’s input. For PPEPs that are longer than the calculated number, we trim the non-template parts that have less importance in the final performance of the model – the parts that do not provide category, score, or title. We select $\bar{L} = 256$ in this paper.

E Performance of the Models on the Validation Set

This section reports the results of experiments on the validation set. Table 5 reports the results of fine-tuning the language model on the user-based separation setting on the validation set. Table 7 shows the results of fine-tuning the language model on the time-based separation setting on the validation set. Table 6 shows the results of zero-shot evaluation of large language models on the user-based separation setting on the validation set. Table 8 depicts the results of zero-shot evaluation of large language models on the time-based separation setting on the validation set.

F Statistics of the Tasks in the LaMP Benchmark

This section reports the statistics of datasets for each task in LaMP. The statistics are reported in Table 9.

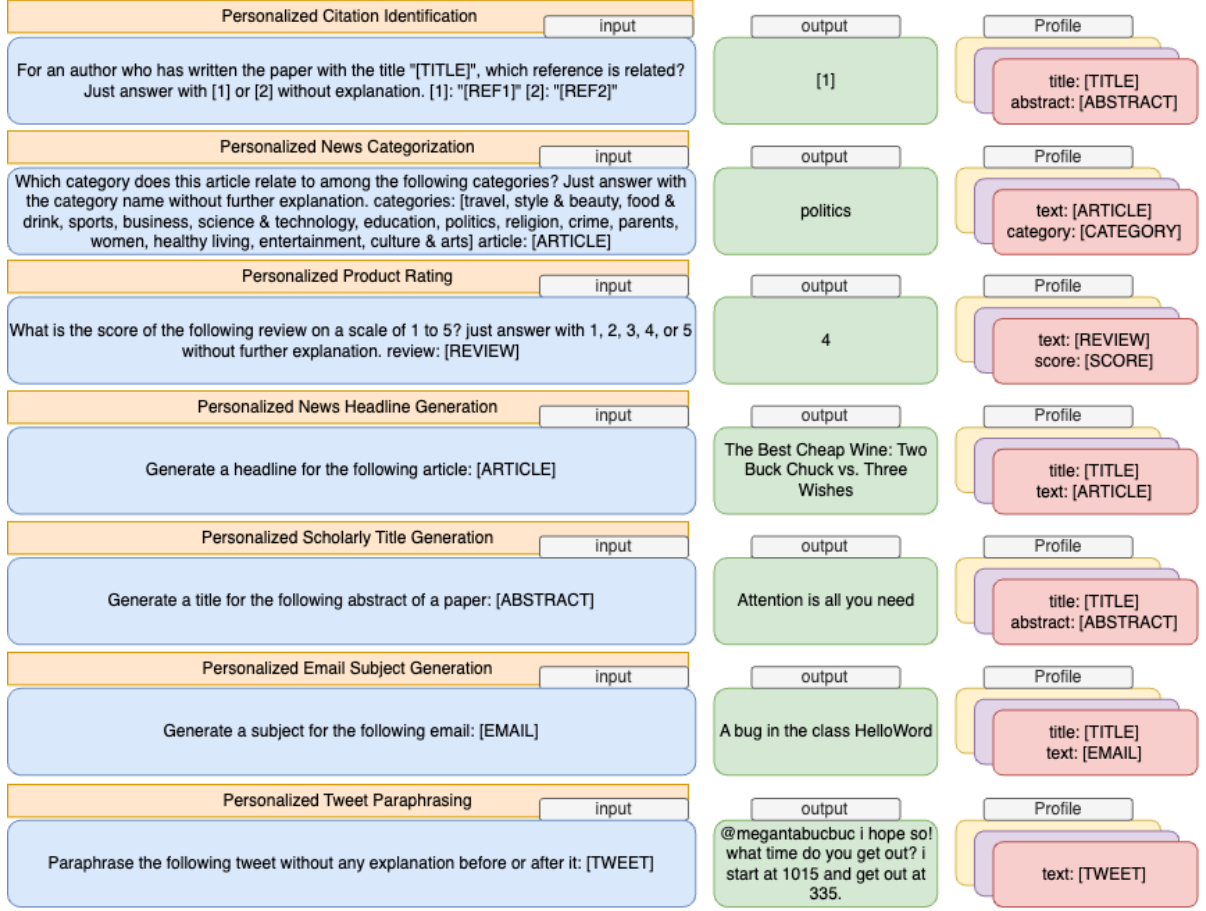


Figure 3: An overview of the templates used for creating data samples for each task in LaMP.

Task	Per Profile Entry Prompt (PPEP)	Aggregated Input Prompt (AIP)
LaMP-1: Citation Ident.	" P_i [title]"	<code>add_to_paper_title(concat([PPEP(P_1), ..., PPEP(P_n)], ", and "), [INPUT])</code>
LaMP-2: News Cat.	the category for the article: " P_i [text]" is " P_i [category]"	<code>concat([PPEP(P_1), ..., PPEP(P_n)], ", and "). [INPUT]</code>
LaMP-3: Product Rat.	P_i [score] is the score for " P_i [text]"	<code>concat([PPEP(P_1), ..., PPEP(P_n)], ", and "). [INPUT]</code>
LaMP-4: News Headline	" P_i [title]" is the title for " P_i [text]"	<code>concat([PPEP(P_1), ..., PPEP(P_n)], ", and "). [INPUT]</code>
LaMP-5: Scholarly Title	" P_i [title]" is the title for " P_i [abstract]"	<code>concat([PPEP(P_1), ..., PPEP(P_n)], ", and "). Following the given patterns [INPUT]</code>
LaMP-6: Email Subject	" P_i [title]" is the title for " P_i [text]"	<code>concat([PPEP(P_1), ..., PPEP(P_n)], ", and "). [INPUT]</code>
LaMP-7: Tweet Para.	" P_i [text]"	<code>concat([PPEP(P_1), ..., PPEP(P_n)], ", and ") are written by a person. Following the given patterns [INPUT]</code>

Table 4: Prompts template used to augment the input of the LM with the user profile. `concat` is a function that concatenates the strings in its first argument by placing the string in the second argument between them. `add_to_paper_title` is a function designed to add the string in its first argument to the paper’s title in the Personalized Citation Identification task. `PPEP` is a function that create the prompt for each entry in the retrieved profile entries. `[INPUT]` is the task’s input.

Dataset	Metric	FlanT5-base (fine-tuned)					
		Non-Personalized	Untuned profile, $k = 1$			Tuned retriever, k	Tuned profile
			Random	BM25	Contriever		
LaMP-1U: Personalized Citation Identification	Accuracy	0.522	0.597	0.623	0.695	Contriever, 4	0.731
LaMP-2U: Personalized News Categorization	Accuracy	0.730	0.771	0.784	0.811	Contriever, 4	0.835
	F1	0.504	0.529	0.555	0.595		0.637
LaMP-3U: Personalized Product Rating	MAE	0.314	0.312	0.282	0.275	Contriever, 4	0.258
	RMSE	0.624	0.633	0.609	0.589		0.572
LaMP-4U: Personalized News Headline Generation	ROUGE-1	0.158	0.167	0.176	0.188	Contriever, 4	0.201
	ROUGE-L	0.144	0.152	0.161	0.172		0.185
LaMP-5U: Personalized Scholarly Title Generation	ROUGE-1	0.424	0.389	0.441	0.405	BM25, 2	0.453
	ROUGE-L	0.382	0.352	0.401	0.367		0.414
LaMP-6U: Personalized Email Subject Generation	ROUGE-1	0.392	0.469	0.575	0.567	BM25, 4	0.583
	ROUGE-L	0.374	0.454	0.563	0.553		0.570
LaMP-7U: Personalized Tweet Paraphrasing	ROUGE-1	0.511	0.512	0.520	0.522	Contriever, 4	0.526
	ROUGE-L	0.456	0.457	0.465	0.467		0.471

Table 5: The personalized text classification and generation results for a fine-tuned language model (i.e., FlanT5-base) on the validation set of user-based separation setting. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3U dataset. k denotes the number of documents retrieved for personalizing language model outputs.

Dataset	Metric	Non-Personalized		Personalized	
		FlanT5-XXL	GPT-3.5	FlanT5-XXL	GPT-3.5
LaMP-1U: Personalized Citation Identification	Accuracy	0.522	0.510	0.675	0.701
LaMP-2U: Personalized News Categorization	Accuracy	0.591	0.610	0.598	0.693
	F1	0.463	0.455	0.471	0.455
LaMP-3U: Personalized Product Rating	MAE	0.357	0.699	0.282	0.658
	RMSE	0.666	0.977	0.5841	1.102
LaMP-4U: Personalized News Headline Generation	ROUGE-1	0.164	0.133	0.192	0.160
	ROUGE-L	0.149	0.118	0.178	0.142
LaMP-5U: Personalized Scholarly Title Generation	ROUGE-1	0.455	0.395	0.467	0.398
	ROUGE-L	0.410	0.334	0.424	0.336
LaMP-6U: Personalized Email Subject Generation	ROUGE-1	0.332	-	0.466	-
	ROUGE-L	0.320	-	0.453	-
LaMP-7U: Personalized Tweet Paraphrasing	ROUGE-1	0.459	0.396	0.448	0.391
	ROUGE-L	0.404	0.337	0.396	0.324

Table 6: The zero-shot personalized text classification and generation results on the validation set of user-based separation setting. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3U dataset. For personalized models, the tuned retriever based on the validation performance was selected.

Dataset	Metric	FlanT5-base (fine-tuned)						
		Non-Personalized	Untuned profile, $k = 1$				Tuned retriever, k	Tuned profile
			Random	BM25	Contriever	Recency		
LaMP-1T: Personalized Citation Identification	Accuracy	0.629	0.662	0.695	0.709	0.681	Contriever, 4	0.732
LaMP-2T: Personalized News Categorization	Accuracy	0.755	0.781	0.788	0.800	0.794	Contriever, 4	0.806
	F1	0.573	0.607	0.621	0.637	0.626		0.647
LaMP-3T: Personalized Product Rating	MAE	0.278	0.273	0.269	0.272	0.260	Recency, 4	0.259
	RMSE	0.595	0.590	0.583	0.589	0.576		0.568
LaMP-4T: Personalized News Headline Generation	ROUGE-1	0.164	0.176	0.176	0.177	0.179	Recency, 2	0.185
	ROUGE-L	0.149	0.160	0.161	0.163	0.165		0.169
LaMP-5T: Personalized Scholarly Title Generation	ROUGE-1	0.462	0.459	0.473	0.470	0.462	Contriever, 4	0.472
	ROUGE-L	0.414	0.412	0.425	0.423	0.416		0.423
LaMP-6T: Personalized Email Subject Generation	ROUGE-1	0.470	0.504	0.509	0.519	0.510	Contriever, 4	0.520
	ROUGE-L	0.455	0.489	0.496	0.507	0.497		0.509
LaMP-7T: Personalized Tweet Paraphrasing	ROUGE-1	0.462	0.507	0.509	0.514	0.510	Contriever, 4	0.518
	ROUGE-L	0.414	0.457	0.460	0.464	0.459		0.467

Table 7: The personalized text classification and generation results for a fine-tuned language model (i.e., FlanT5-base) on the validation set of time-based separation setting. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3U dataset. k denotes the number of documents retrieved for personalizing language model outputs.

Dataset	Metric	Non-Personalized		Personalized	
		FlanT5-XXL	GPT-3.5	FlanT5-XXL	GPT-3.5
LaMP-1T: Personalized Citation Identification	Accuracy	0.498	0.478	0.656	0.640
LaMP-2T: Personalized News Categorization	Accuracy	0.636	0.687	0.587	0.678
	F1	0.541	0.551	0.507	0.556
LaMP-3T: Personalized Product Rating	MAE	0.335	0.720	0.294	0.608
	RMSE	0.639	1.000	0.586	1.022
LaMP-4T: Personalized News Headline Generation	ROUGE-1	0.173	0.146	0.192	0.159
	ROUGE-L	0.157	0.128	0.175	0.138
LaMP-5T: Personalized Scholarly Title Generation	ROUGE-1	0.472	0.413	0.472	0.421
	ROUGE-L	0.419	0.348	0.422	0.352
LaMP-6T: Personalized Email Subject Generation	ROUGE-1	0.316	-	0.382	-
	ROUGE-L	0.302	-	0.369	-
LaMP-7T: Personalized Tweet Paraphrasing	ROUGE-1	0.454	0.390	0.440	0.392
	ROUGE-L	0.401	0.331	0.391	0.325

Table 8: The zero-shot personalized text classification and generation results on the validation set of time-based separation setting. For all metrics the higher the better, except for RMSE and MAE which are used for the LaMP-3T dataset. For personalized models, the tuned retriever based on the validation performance was selected.

Task	Type	Separation	#train	#dev	#test	Input Length	Output Length	#Profiles	#classes
Citation Ident.	binary classification	user	9682	2500	2500	51.40 \pm 5.72	-	90.61 \pm 53.87	2
		time	6542	1500	1500	51.43 \pm 5.70		84.15 \pm 47.54	
News Cat.	categorical classification	user	5914	1052	1274	65.93 \pm 12.29	-	306.42 \pm 286.65	15
		time	8090	1605	1568	65.44 \pm 11.40		191.03 \pm 168.43	
Product Rat.	ordinal classification	user	20000	2500	2500	145.14 \pm 157.96	-	188.10 \pm 129.42	5
		time	20000	2500	2500	128.18 \pm 146.25		185.40 \pm 129.30	
News Headline	text generation	user	12527	1925	2376	30.53 \pm 12.67	9.78 \pm 3.10	287.16 \pm 360.62	-
		time	12500	1500	1800	29.97 \pm 12.09	10.07 \pm 3.10	204.59 \pm 250.75	
Scholarly Title	text generation	user	9682	2500	2500	152.81 \pm 86.60	9.26 \pm 3.13	89.61 \pm 53.87	-
		time	14682	1500	1500	162.34 \pm 65.63	9.71 \pm 3.21	87.88 \pm 53.63	
Email Subject	text generation	user	4840	1353	1246	436.15 \pm 805.54	7.34 \pm 2.83	80.72 \pm 51.73	-
		time	4821	1250	1250	454.87 \pm 889.41	7.37 \pm 2.78	55.67 \pm 36.32	
Tweet Para.	text generation	user	10437	1500	1496	29.76 \pm 6.94	16.93 \pm 5.65	17.74 \pm 15.10	-
		time	13437	1498	1500	29.72 \pm 7.01	16.96 \pm 5.67	15.71 \pm 14.86	

Table 9: Statistics of the datasets in the LaMP benchmark.