# Your Next Token Prediction: A Multilingual Benchmark for Personalized Response Generation

Shiyao Ding
Kyoto University
Kyoto, Japan
ding@i.kyoto-u.ac.jp

Takayuki Ito
Kyoto University
Kyoto, Japan
ito@i.kyoto-u.ac.jp

## ABSTRACT

Large language models (LLMs) excel at general *next-token prediction* but still struggle to generate responses that reflect how individuals truly communicate, such as replying to emails or social messages in their own style. However, real SNS or email histories are difficult to collect due to privacy concerns. To address this, we propose the task of **"Your Next Token Prediction (YNTP)"**, which models a user's precise word choices through controlled human–agent conversations. We build a multilingual benchmark of **100 dialogue sessions** across English, Japanese, and Chinese, where users interact for five days with psychologically grounded NPCs based on MBTI dimensions. This setup captures natural, daily-life communication patterns and enables analysis of users' internal models. We evaluate prompt-based and fine-tuning-based personalization methods, establishing the first benchmark for YNTP and a foundation for user-aligned language modeling. The dataset is available at: https://github.com/AnonymousHub4Submissions/your-next-token-prediction-dataset-100.

## KEYWORDS

Large language models, Personalized alignment, Human-agent dialogue

## 1 INTRODUCTION

Large language models (LLMs) are trained on the objective of *next-token prediction*, a formulation that has enabled remarkable general capabilities in text generation, reasoning, and dialogue. Yet despite their fluency, most LLMs remain *impersonal*—they produce contextually appropriate but socially generic responses that ignore the distinctive preferences, tones, and values of individual users. As a result, when applied to personalized communication (e.g., daily chat, email, or social messaging), LLMs tend to generate replies that "sound like a model," not like the user.

Recent progress in **personalized alignment** seeks to bridge this gap by conditioning generation on user identity or preference data. Existing benchmarks such as **LaMP/LongLaMP** [8, 14], **P-SOUPS** [6], and **PRISM** [7] have advanced evaluation of personalization across domains including news, reviews, and dialogue. However, these datasets largely focus on English text, short or task-oriented inputs, and surface-level style imitation (e.g., headline rewriting or response rephrasing). They do not capture the dynamic, emotionally grounded conversations that characterize human interactions, nor do they model the *token-level decisions* that reflect each user's internal cognitive and affective patterns.

Thus, our goal is to make LLMs capable of replying to messages or emails in the same way a specific user would, by understanding the user's internal model. However, real SNS or email histories that reflect genuine personal communication are extremely difficult to obtain due to privacy and ethical constraints. To overcome this limitation, we construct a controlled **human–agent conversation environment**, where users engage in natural daily dialogues with psychologically grounded agents that simulate interpersonal interactions.

Specifically, we introduce the first benchmark for **"your next-token prediction"**—a fine-grained formulation of personalized language modeling where the goal is to predict how an individual would respond, word by word, to a given message. Our benchmark is constructed from **multilingual human–agent conversations** with 100 participants (English, Japanese, and Chinese, 30+ each) engaging in **five-day dialogue sessions** with LLM-driven non-player characters (NPCs). Unlike prior datasets built from static corpora, our dialogues simulate **daily-life interactions**—similar to SNS or chat-app exchanges—with NPCs designed around distinct personalities and psychological roles.

Each NPC follows a structured finite-state machine (FSM) and is grounded in the **MBTI framework**, allowing targeted exploration of user traits such as extraversion or intuition through natural conversation. This interactive design enables extraction of each user's *internal model*—their consistent linguistic, emotional, and decision-making patterns—across multiple days and contexts. The benchmark task is to predict a user's fifth-day response given their preceding history, providing a testbed for models to capture continuity, adaptation, and personal consistency.

We evaluate both **prompt-based** and **fine-tuning-based** alignment methods on this dataset, establishing the first quantitative baseline for personalized response generation at the token level. By combining multilingual daily dialogue, psychologically grounded design, and multi-day interaction, this benchmark moves beyond stylistic mimicry toward modeling the deeper cognitive regularities that govern how individuals choose words. Ultimately, *your next-token prediction* represents a critical step toward the "last mile" of alignment—transforming LLMs from generic communicators into authentic personal agents. The dataset is available at: https://github.com/AnonymousHub4Submissions/your-next-token-prediction-dataset-100.

## 2 RELATED WORK

Research on personalized alignment in natural language processing has developed along two major lines: datasets that enable systematic evaluation of such methods, and methods for adapting large language models (LLMs) to user-specific preferences. In this section, we summarize both aspects.

### 2.1 Personalized Alignment Datasets

Datasets play a central role in benchmarking personalization. Several corpora have been introduced to capture user-level text and enable systematic evaluation. We organize them into two groups: (1) *personalized benchmarks* that explicitly condition generation on user identity or style, and (2) *preference-based datasets* that use human feedback for alignment at the population level. A summary is shown in Table 1.

**Personalized Benchmarks.** The most established family of datasets is the **LaMP Benchmarks** [14], which unify seven personalization tasks such as personalized title generation and email subject rewriting, all conditioned on user-specific histories. **LongLaMP** [8] extends this benchmark to long-form generation tasks such as reviews and blog posts, evaluated by BLEU, ROUGE, and METEOR. These benchmarks standardized evaluation and highlighted the importance of retrieval and user history usage, though they primarily demonstrate how users interact with content rather than capturing deeper internal preferences. Beyond LaMP, **P-SOUPS** [6] expands personalization across multiple domains, integrating user profiles and style history to evaluate multi-task alignment. **PRISM** [7] focuses on next-turn personalized response generation, constructing large-scale dialogue data annotated with user-level stylistic metadata. **PersonalLLM** [20] and **PERSONA** [1] further explore user embeddings and persona-consistent generation, representing recent advances in explicit user modeling for generative LLMs. Together, these datasets form the core of personalized alignment research by modeling how LLMs adapt to user-specific writing styles and behavioral traits.

**Preference-Based and Feedback Datasets.** Several corpora investigate alignment through human or model preference signals rather than user-specific identities. **FLASK** [18], **REGEN** [15], and **ALOE** [17] collect human preference annotations or evolving dialogue feedback to improve alignment quality. **PREFEVAL** [19] provides standardized evaluation for model preference judgments. While such datasets contribute to alignment research, they primarily address population-level preference optimization (e.g., RLHF or P-RLHF) rather than personalized response modeling targeted in this work.

**Comparison with Our Benchmark.** Unlike previous datasets such as LaMP, P-SOUPS, and PRISM—which focus on English, single-turn, or task-specific personalization—our **Your Next Token Prediction (YNTP)** benchmark captures *multi-day, multilingual human–agent conversations* grounded in daily-life contexts. It integrates psychological modeling (MBTI dimensions) within natural dialogues to elicit each user's internal traits, enabling token-level prediction of individualized responses. This design moves beyond surface style imitation toward modeling deeper cognitive and behavioral consistency in personalized language generation.

### 2.2 Personalized Alignment Methods

Personalized alignment methods adapt large language models (LLMs) to user preferences and histories, enabling contextually consistent and user-specific outputs. These approaches fall into two paradigms: **(1) external personalization**, which conditions generation without altering model parameters, and **(2) internal personalization**, which updates parameters for individualized behavior.

*(1) External Personalization (Without Modifying Parameters).* External approaches treat the LLM as frozen and inject user information through inputs or inference-time control. **Prompt-based methods** directly incorporate user histories or summaries into prompts [3, 12, 16]. Retrieval-Augmented Generation (RAG) extends this by retrieving user-specific documents or prior responses as in-context demonstrations [14], sometimes optimizing retrieval with generation-based rewards [13]. Another direction employs **persona-conditioned inference**, representing each user with embeddings or structured profiles such as Persona-Plug [10]. Latent steering and decoding-time editing further adjust internal activations to reflect user traits—achieving scalable personalization without retraining or privacy risks.

*(2) Internal Personalization (With Modifying Parameters).* Internal methods directly adapt LLM parameters to encode user-specific patterns. **Parameter-efficient fine-tuning (PEFT)** approaches such as LoRA and prefix-tuning [4, 5] update a small subset of weights for user adaptation while limiting computation. **Personalized RLHF (P-RLHF)** [9] further tunes reward models from user feedback, enabling fine-grained alignment but requiring costly data collection. Recent inference-stage variants—e.g., PAD (Personalized Alignment at Decoding) [2] and CHAMELEON [11]—steer token probabilities or latent activations via personalized rewards, achieving lightweight but effective adaptation.

Overall, external methods emphasize flexibility and privacy, while internal methods enable deeper, parameter-level alignment. Together, they outline complementary pathways toward scalable personalized language modeling.

## 3 PROBLEM DEFINITION: YOUR NEXT TOKEN PREDICTION

We define the personalized response generation task as **Your Next Token Prediction**—given an incoming message from others, the goal is to predict how a specific user would respond, in both content and style.

*Task Formulation.* Let $x$ denote an incoming message or dialogue context, and $y^u = (y_1^u, \ldots, y_m^u)$ represent the actual response written by user $u$. The model aims to generate a personalized response $\hat{y}^u = (\hat{y}_1^u, \ldots, \hat{y}_m^u)$ conditioned on both the input message and the user's background information $U^u$, which includes profile features $p^u$ (e.g., personality, preference) and historical dialogue $H^u$:

$$U^u = (p^u, H^u).$$

The probabilistic objective of user-conditioned generation is thus:

$$P_\theta(\hat{y}^u \mid x, U^u) = \prod_{t=1}^{m} P_\theta(\hat{y}_t^u \mid \hat{y}_{<t}^u, x, U^u),$$

**Table 1: Summary of key personalization-related datasets.**

| Dataset | Input | Output | History/Persona | Size (Train/Dev) | Metric |
|---|---|---|---|---|---|
| LaMP (1–7) | Multi-domain prompts (papers, movies, reviews, news, tweets) | Personalized or preference-aware outputs (titles, tags, ratings, headlines, paraphrases) | User-specific histories (authored papers, rated items, past headlines, tweets, etc.) | 6k–20k per task | Accuracy, F1, MAE, ROUGE-1/L |
| LongLaMP | Long-form texts (reviews, blogs) | Personalized long-form outputs | User's prior writings | ~10k examples | BLEU, ROUGE, METEOR |
| P-SOUPS | Prompts across news/dialogue/review domains | Personalized responses | User profiles + multi-task history | ~1M examples (2024) | BLEU, ROUGE, human eval. |
| PRISM | Conversational context (chat logs) | Next-turn personalized response | User style embeddings + history | ~50k conversations | BLEU, StyleSim, human eval. |
| CUSTOM | Email prompts | Personalized email response | Few-shot (4 train, 1 test per user) | 2 prompts × 2 demos + 1 test per user | Human preference, style match |
| PersonalLLM | Generic task prompt | Personalized response | Learned user embeddings | ~100k samples | BLEU, style match, human eval. |
| PERSONA | Task prompt + persona card | Persona-consistent response | Explicit persona description | ~60k examples | StyleSim, human eval. |
| FLASK | Task + user preference | Aligned model response | Human-annotated preferences | ~90k examples | BLEU, ROUGE, human eval. |
| REGEN | Dialogue context | Personalized continuation | Past dialogues with evolving preferences | ~50k dialogues | Consistency, human eval. |
| ALOE | Prompt + rationale | Personalized response + explanation | User-specific rationales | ~80k examples | BLEU, ROUGE, faithfulness |
| PREFEVAL | Candidate generations | Human preference ranking | Annotator preference judgments | ~20k comparisons | Win-rate, agreement |

where $\theta$ denotes the parameters of the large language model (LLM). This formulation extends standard next-token prediction by explicitly conditioning on the user dimension $u$, enabling the model to produce responses that align with how the specific user would phrase their next utterance.

*Learning Objective.* Given a dataset of user-specific dialogue pairs

$$\mathcal{D} = \{(x_i, y_i^u, U_i^u)\}_{i=1}^N,$$

the training objective minimizes the discrepancy between the predicted response $\hat{y}_i^u$ and the user's real response $y_i^u$. The loss function follows a token-level negative log-likelihood (cross-entropy) form:

$$\min_\theta \ \mathcal{L}_{\text{YNP}}(\theta) = -\mathbb{E}_{(x,y^u,U^u)\sim\mathcal{D}} \left[ \sum_{t=1}^{|y^u|} \log P_\theta(y_t^u \mid y_{<t}^u, x, U^u) \right].$$

This objective encourages the model to generate responses whose token distributions are close to those of the real user, effectively minimizing the divergence between predicted and true personalized responses.
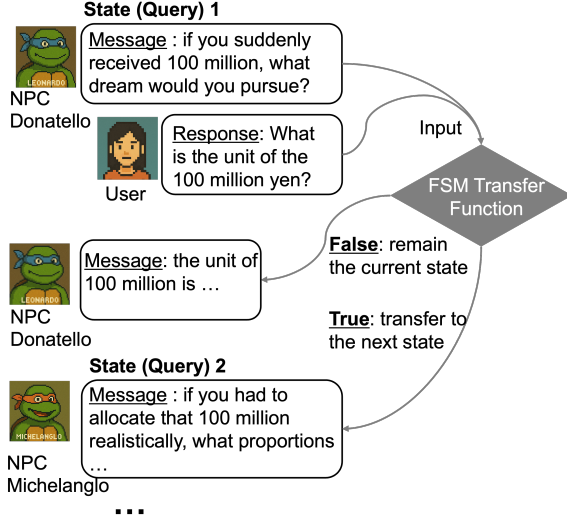
*Interpretation.* Under this formulation, **Your Next Token Prediction** generalizes classical next-token prediction to a personalized setting, where each user's communication history and profile shape the conditional probability of every token. The model thereby learns not only *what* the user tends to say (semantic preference) but also *how* they express it (stylistic preference).

## 4 LLM-DRIVEN MULTI-NPC DIALOGUE SYSTEM FOR BENCHMARK CONSTRUCTION

To realize the goal of constructing a benchmark for **Your Next Token Prediction (YNP)**, we develop an **LLM-driven multi-NPC dialogue system** that automatically collects natural, personalized, and psychologically grounded conversation data. The system simulates a multi-day shared-house experience where users interact with multiple non-player characters (NPCs), each driven by large language models (LLMs) and coordinated through a finite state machine (FSM). Through these structured yet dynamic interactions, we collect paired message–response data $\{(x, y^u, U^u)\}$ necessary for training and evaluating personalized alignment methods.

*Personality Modeling via MBTI..* To systematically capture user individuality, the system adopts the **MBTI** framework as an interpretable representation of users' internal models. Each MBTI dimension (E/I, S/N, T/F, J/P) corresponds to distinct reasoning or affective patterns that shape communication styles. **Figure 1** illustrates the FSM control flow: each state is a query probing one MBTI dimension. A user's response is evaluated by a check function—if it sufficiently reflects the targeted trait, the dialogue proceeds to the next state; otherwise, the NPC triggers reflective feedback to prompt elaboration.

*Framework Overview.* The system integrates three coordinated components:

**Figure 1: FSM-based control flow of the LLM-driven multi-NPC dialogue system. Each state corresponds to a query probing one MBTI dimension. User responses are evaluated by a check function that determines transition or reflection.**



**Figure 2: Dialogue interface (shared-house environment and status panel) visualizing user–NPC interactions across multiple personality dimensions.**

- **FSM Engine:** A Python-based FSM governs dialogue flow and state transitions across NPCs. Each NPC maintains an individualized path reflecting personality traits and MBTI dimensions. The FSM includes a **transfer function** that evaluates user responses: if the response aligns with the targeted MBTI trait, the dialogue advances to the next state (*True*); otherwise, it remains in the current state (*False*) to elicit clarification or reflection.
- **Scenario Script:** A JSON-based scenario file defines the dialogue content, branching logic, and NPC roles that shape the overall narrative flow. Each state specifies the NPC speaker, targeted MBTI dimension, and possible user response patterns. By modifying this script, researchers can easily construct new domains (e.g., workplace communication, counseling, education) or customize interaction depth, number of days, and NPC personalities without altering system code. The script therefore serves as a flexible interface between experimental design and the underlying FSM engine.
- **LLM Dialogue Generation:** The LLM layer functions as the linguistic and emotional realization module on top of the FSM logic. Given the current state, user input, and dialogue history, the model dynamically adjusts tone, phrasing, and elaboration to reflect contextual nuances and inferred user traits. For instance, when interacting with an introverted user, the LLM may adopt a gentler tone and shorter responses, whereas with an extroverted user it may generate more expressive and expansive utterances. This enables contextually adaptive NPC behaviors that go beyond static scripts and produce natural, human-like multi-turn conversations.

*Example Scenario.* As shown in Figure 1 and surfaced through the interface in Figure 2, when the NPC *Donatello* asks, "If you suddenly received 100 million yen, what dream would you pursue?", a user might reply "I want to travel." NPCs react in character: *Donatello* emphasizes financial planning, while *Michelangelo* highlights adventure—creating realistic, multi-perspective conversations.

*Data Collection for Benchmarking.* Each user–NPC interaction yields a pair $(x, y^u)$ with metadata $U^u$ (day index, NPC identity, MBTI state). The dataset accumulates as

$$\mathcal{D}^u = \{(x_t, y_t^u, U_t^u)\}_{t=1}^{T_u},$$

supporting the YNP objective:

$$\min_\theta \mathcal{L}_{\mathrm{YNP}}(\theta) = -\mathbb{E}_{(x,y^u,U^u)\sim\mathcal{D}}\left[\sum_{t=1}^{|y^u|} \log P_\theta(y_t^u \mid y_{<t}^u, x, U^u)\right].$$

*User Data Structure.* All dialogues are stored in a structured JSON format, exemplified in **Figure 3**. Each record contains metadata (game_number, start_time, completion_time, language, status) and daily logs (day_1, day_2, ...) with ordered message–response pairs and speaker identities. For benchmark construction, we use the dialogues from the **first four days** of each user as the **training set** to model their linguistic and behavioral patterns, and reserve the **fifth day** as the **test set**. During evaluation, the model generates predicted responses to fifth-day messages, which are then compared with the real user responses to measure alignment accuracy and stylistic consistency.

## 5 EXPERIMENTS

In this section, we conduct extensive experiments across English, Chinese, and Japanese user groups to evaluate the effectiveness of our personalized response generation framework. Using multiple large language models (LLMs) and prompting strategies, we assess both semantic fidelity and stylistic alignment under the proposed 2S evaluation framework.

**Table 2: Evaluation metrics for personalized response generation.**

| Group | ID | Metric Name | Description |
|---|---|---|---|
| **Substance (What to say)** | M1 | Word Mover's Distance | Measures semantic dissimilarity between predicted and ground-truth responses based on the minimal cumulative distance of word embeddings. Lower values indicate higher semantic overlap. |
| | M2 | Sentence Similarity | Cosine similarity between sentence embeddings of generated and real responses, capturing overall semantic closeness and contextual coherence. |
| | M3 | Content Similarity | GPT-5–assessed similarity focusing on factual and conceptual consistency between generated and ground-truth responses. |
| **Style (How to say)** | M4 | Normalized Length Similarity | Compares response lengths after normalization to evaluate verbosity or conciseness consistency between users and models. |
| | M5 | Style Similarity | GPT-5–based stylistic evaluation including tone, formality, sentiment, and expressiveness alignment between user and model responses. |
| | M6 | History Similarity | Cosine similarity between embeddings of generated responses and users' previous responses, measuring continuity and self-consistency. |

## 5.1 Experiment Settings

In this subsection, we describe the detailed experimental setup of our benchmark, including the evaluated models, datasets, baseline methods, and evaluation metrics.

**Evaluated Models and Datasets.** We evaluate a diverse set of large language models (LLMs) and prompting configurations across three languages—English, Chinese, and Japanese—to ensure generalizability and cross-lingual robustness. The evaluated models include `gpt-3.5-turbo`, `gpt-4o-mini`, `gpt-4.1-mini`, and `gpt-5-mini`, representing a progression from mid-scale to cutting-edge reasoning capabilities. In addition, we include open-source models such as `Claude`, `Gemini`, and a parameter-efficient fine-tuning model `L3-ELYZA-8B(LoRA)` to assess the adaptability of both commercial and open architectures.

To provide a comprehensive evaluation, we construct three multilingual datasets derived from our **multi-day, multi-NPC dialogue system**. Each dataset corresponds to one language and contains user–NPC interaction logs collected over five simulated days, where users converse with multiple NPCs under a shared-house scenario. Every conversation turn is stored in JSON format, including NPC messages, user replies, timestamps, and dialogue states. In total, the benchmark comprises 34 English users, 33 Chinese users, and 33 Japanese users, ensuring balanced linguistic and cultural representation. These datasets are used purely for evaluation without any model fine-tuning.

**Baseline Methods.** We compare our proposed framework with three baseline prompting strategies to analyze the contribution of few-shot examples and internal-model reasoning:

- **Prompt Engineering (zero-shot):** The LLM directly generates personalized responses without any examples, relying solely on the system and user prompts.
- **Prompt Engineering (few-shot):** The LLM receives several demonstration examples drawn from previous user–NPC interactions to enhance personalization consistency.
- **Prompt Engineering (few-shot with MBTI inference):** The LLM is augmented with explicit reasoning over the user's inferred MBTI dimensions, enabling adaptive personality-aware responses.

- **PEFT (LoRA):** We further fine-tune an LLM using personalized data, designed to test whether lightweight adaptation provides additional stylistic alignment advantages beyond in-context prompting.

*Evaluation Metrics.* We evaluate personalized response generation under the guiding **2S Principle**—*Substance* (what to say) and *Style* (how to say). These two complementary dimensions jointly assess whether model outputs are both semantically faithful to the intended meaning and stylistically aligned with the target user. Table 2 summarizes the six evaluation metrics used in our benchmark.

*Substance metrics (M1–M3)* capture semantic and factual fidelity. **M1 (Word Mover's Distance)** measures the minimal transport distance between word embeddings of generated and gold responses. **M2 (Sentence Similarity)** evaluates contextual coherence via cosine similarity of sentence embeddings. **M3 (Content Similarity)** uses GPT-5 as a semantic judge to evaluate factual and conceptual alignment.

*Style metrics (M4–M6)* quantify stylistic and behavioral consistency. **M4 (Normalized Length Similarity)** assesses verbosity/conciseness alignment. **M5 (Style Similarity)** measures tone and formality matching through GPT-5–based judgments. **M6 (History Similarity)** evaluates whether generated responses remain consistent with a user's previous communication history.

In addition to automatic metrics, both **human** and **GPT-5–based** judges assess outputs under the same 2S framework. They rate **Substance**—semantic appropriateness and factual accuracy—and **Style**—whether the reply "sounds like" the same user—on a 1–5 Likert scale.

To compare models, we also report a **pairwise win rate** metric:

$$\text{WinRate} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}\big(\text{pref}(\hat{y}_i^A, \hat{y}_i^B) = A\big),$$

which reflects the proportion of cases in which model $A$'s output is preferred over model $B$ by human or GPT-5 evaluators.

Table 3: Main results across three languages: macro-averaged scores over 33 English users, 34 Chinese users, and 33 Japanese users. Each cell represents the average score of all participants in that language group on the corresponding metric. Lower values of M1 indicate better distance alignment, while higher values of M2–M6 indicate stronger similarity or consistency.

| Method | Base Model | English (33 users) | | | | | | Chinese (34 users) | | | | | | Japanese (33 users) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | M1↓ | M2↑ | M3↑ | M4↑ | M5↑ | M6↑ | M1↓ | M2↑ | M3↑ | M4↑ | M5↑ | M6↑ | M1↓ | M2↑ | M3↑ | M4↑ | M5↑ | M6↑ |
| Prompt Eng. (zero-shot) | gpt-3.5-turbo | 0.242 | 0.427 | 3.092 | 0.337 | 2.876 | 0.430 | 1.042 | 0.417 | 3.052 | 0.252 | 2.192 | 0.461 | 0.999 | 0.178 | 3.105 | 0.252 | 1.983 | 0.289 |
| | gpt-4o-mini | 0.242 | 0.425 | 3.152 | 0.236 | 2.837 | 0.412 | 0.928 | 0.538 | 3.179 | 0.226 | 2.497 | 0.723 | 0.717 | 0.277 | 3.252 | 0.236 | 2.572 | 0.582 |
| | gpt-4.1-mini | 0.249 | 0.427 | 3.180 | 0.266 | 2.839 | 0.413 | 0.944 | 0.524 | 3.164 | 0.221 | 2.578 | 0.712 | 0.721 | 0.310 | 3.318 | 0.246 | 2.897 | 0.594 |
| | gpt-5-mini | 0.264 | 0.373 | 3.049 | 0.082 | 2.373 | 0.405 | 1.042 | 0.376 | 3.197 | 0.050 | 2.051 | 0.597 | 0.836 | 0.175 | **3.325** | 0.064 | 2.622 | 0.443 |
| | gemini-2.5-flash | 0.244 | 0.389 | 3.243 | 0.121 | 2.171 | 0.426 | 0.983 | 0.406 | **3.289** | 0.072 | 1.964 | 0.631 | 0.751 | 0.188 | 3.196 | 0.086 | 1.893 | 0.483 |
| | claude-sonnet-4-0 | **0.235** | 0.367 | 3.176 | 0.105 | 1.821 | 0.436 | 1.043 | 0.351 | 3.242 | 0.098 | 1.673 | 0.552 | 0.812 | 0.177 | 3.211 | 0.131 | 2.016 | 0.470 |
| Prompt Eng. (few-shot) | gpt-3.5-turbo | 0.249 | 0.437 | 3.087 | 0.382 | 2.826 | 0.445 | 0.921 | 0.606 | 3.189 | 0.340 | 2.894 | 0.753 | 0.717 | 0.450 | 3.174 | 0.367 | 3.011 | 0.624 |
| | gpt-4o-mini | 0.247 | 0.451 | 3.204 | 0.368 | 3.274 | 0.428 | 0.908 | 0.615 | 3.189 | 0.327 | 3.161 | 0.745 | 0.688 | 0.464 | 3.198 | 0.349 | 4.225 | 0.657 |
| | gpt-4.1-mini | 0.243 | **0.456** | **3.288** | 0.335 | 3.279 | 0.414 | 0.929 | 0.610 | 3.214 | 0.303 | 3.122 | 0.758 | 0.700 | 0.422 | 3.287 | 0.312 | 3.994 | 0.637 |
| | gpt-5-mini | 0.265 | 0.384 | 2.921 | 0.250 | 3.219 | 0.420 | 0.997 | 0.463 | 3.205 | 0.183 | 2.870 | 0.690 | 0.794 | 0.215 | 3.310 | 0.143 | 3.714 | 0.510 |
| | gemini-2.5-flash | 0.276 | 0.411 | 2.917 | 0.621 | 3.652 | 0.417 | 0.913 | 0.628 | 3.098 | 0.536 | **3.693** | 0.789 | 0.694 | 0.558 | 3.134 | 0.576 | 4.316 | 0.725 |
| | claude-sonnet-4-0 | 0.246 | 0.433 | 3.038 | 0.426 | 3.567 | 0.433 | **0.888** | 0.629 | 3.216 | 0.436 | 3.747 | 0.786 | **0.683** | 0.540 | 3.094 | 0.507 | **4.486** | 0.714 |
| Prompt Eng. (few-shot + MBTI inference) | gpt-3.5-turbo | 0.247 | 0.402 | 2.972 | 0.346 | 2.546 | 0.398 | 0.950 | 0.564 | 3.080 | 0.305 | 2.600 | 0.699 | 0.733 | 0.441 | 3.187 | 0.370 | 3.268 | 0.616 |
| | gpt-4o-mini | 0.240 | 0.438 | 3.206 | 0.347 | 3.184 | 0.425 | 0.919 | 0.615 | 3.199 | 0.293 | 3.012 | 0.748 | 0.689 | 0.434 | 3.212 | 0.324 | 4.020 | 0.639 |
| | gpt-4.1-mini | 0.237 | 0.442 | 3.202 | 0.372 | 3.319 | 0.412 | 0.923 | 0.622 | 3.197 | 0.333 | 3.230 | 0.765 | 0.702 | 0.448 | 3.227 | 0.329 | 4.196 | 0.647 |
| | gpt-5-mini | 0.257 | 0.389 | 3.096 | 0.224 | 3.022 | 0.429 | 0.995 | 0.470 | 3.241 | 0.155 | 2.883 | 0.682 | 0.792 | 0.208 | 3.268 | 0.142 | 3.787 | 0.507 |
| | gemini-2.5-flash | 0.273 | 0.428 | 2.935 | **0.642** | **3.697** | 0.421 | 0.899 | **0.644** | 3.050 | 0.538 | 3.650 | 0.787 | 0.692 | 0.528 | 3.151 | 0.513 | 4.227 | 0.711 |
| | claude-sonnet-4-0 | 0.244 | 0.413 | 3.084 | 0.344 | 3.260 | 0.424 | 0.905 | 0.612 | 3.144 | 0.364 | 3.671 | 0.770 | 0.695 | 0.480 | 3.098 | 0.428 | 4.380 | 0.681 |
| PEFT (LoRA) | L3-ELYZA-8B | 0.273 | 0.349 | 2.323 | 0.635 | 3.493 | **0.450** | 0.942 | 0.641 | 2.585 | **0.561** | 3.541 | **0.791** | 0.737 | **0.571** | 2.584 | **0.621** | 4.137 | **0.760** |

By combining these six quantitative metrics and pairwise preference judgments, our evaluation framework provides a comprehensive and interpretable assessment of personalized response generation across both semantic and stylistic dimensions.

## 5.2 Experiment Results

*5.2.1 Overall Analysis.* All scores are *macro-averaged over participants* within each language group (English: 33, Chinese: 34, Japanese: 33). M1 is a distance metric (*lower is better*); M2–M6 are similarity-style metrics (*higher is better*). We compare four method families: **PE (0)** = Prompt Engineering (zero-shot), **PE (F)** = Prompt Engineering (few-shot), **PE (F+MBTI)** = Prompt Engineering (few-shot + MBTI inference), and **LoRA** = PEFT (LoRA). Base models include gpt-3.5/4o/4.1/5-mini and gemini-2.5-flash, claude-sonnet-4-0, where the results are summarized in Table 3.

*(1) Substance (M1–M3).* Across all three languages, few-shot prompting consistently improves **M2** (sentence similarity) and **M3** (content similarity) over zero-shot, confirming the value of in-context personalization. In English and Chinese, the strongest substance is achieved by **PE (F)** and **PE (F+MBTI)** on top of mid-to-strong bases (gpt-4o/4.1, gemini-2.5-flash, claude-sonnet-4-0), with the best M2 around the mid–0.4s and the best M3 ≈3.2–3.3. Japanese exhibits the same pattern, with gpt-4.1-mini and the two new proprietary models competitive under few-shot. On **M1** (WMD), **PE (F+MBTI)** often attains the lowest values, suggesting that even a coarse personality hypothesis helps the model choose words closer to a user's lexical preferences.

*(2) Style & Consistency (M4–M6).* Stylistic control benefits from two routes: (i) **PE (F)** steadily raises **M4** (length/format consistency) and **M6** (consistency with user history), and (ii) **LoRA** provides strong **M5** (style similarity) and competitive M6, especially in Chinese and Japanese. Among proprietary bases, gemini-2.5-flash and claude-sonnet-4-0 under few-shot usually match or surpass open baselines on M4–M6, with **PE (F+MBTI)** giving an additional boost on M5 when stylistic cues matter.

*(3) Cross-language trends.* English shows the most balanced gains across M2–M6 with few-shot prompting. Chinese benefits markedly on style metrics (M5–M6), indicating that discourse markers and tone cues are well captured by both LoRA and few-shot. Japanese has slightly lower M4 (length normalization), but few-shot methods still dominate zero-shot on all metrics.

*(4) Method takeaways.* **PE (F)** is a strong default for semantic fidelity (M2–M3) and formatting (M4); **PE (F+MBTI)** further helps lexical/style choices (best M1, strong M3/M5); **LoRA** is preferred when sustained style imitation (M5) and historical consistency (M6) are crucial. No single method wins every metric; selection should be *metric- and evidence-driven* (available demos vs. lightweight tuning).

*5.2.2 Heatmap Analysis.* Figure 3 shows average scores across six metrics (M1–M6) for all method–model pairs in English, Chinese, and Japanese. Each cell represents the mean over participants in that language group.

**Performance Heatmap (English)**

| Method-Base Model | M1 | M2 | M3 | M4 | M5 | M6 |
|---|---|---|---|---|---|---|
| LoRA-L3-ELYZA-8B | 0.27 | 0.35 | 2.32 | 0.64 | 3.49 | 0.45 |
| PE (0)-3.5-turbo | 0.24 | 0.43 | 3.09 | 0.34 | 2.88 | 0.43 |
| PE (0)-4.1-mini | 0.25 | 0.43 | 3.18 | 0.27 | 2.84 | 0.41 |
| PE (0)-4o-mini | 0.24 | 0.42 | 3.15 | 0.24 | 2.84 | 0.41 |
| PE (0)-5-mini | 0.26 | 0.37 | 3.05 | 0.08 | 2.37 | 0.41 |
| PE (0)-claude-sonnet-4-0 | 0.23 | 0.37 | 3.18 | 0.10 | 1.82 | 0.44 |
| PE (0)-gemini-2.5-flash | 0.24 | 0.39 | 3.24 | 0.12 | 2.17 | 0.43 |
| PE (F)-3.5-turbo | 0.25 | 0.44 | 3.09 | 0.38 | 2.83 | 0.45 |
| PE (F)-4.1-mini | 0.24 | 0.46 | 3.29 | 0.34 | 3.28 | 0.41 |
| PE (F)-4o-mini | 0.25 | 0.45 | 3.20 | 0.37 | 3.27 | 0.43 |
| PE (F)-5-mini | 0.27 | 0.38 | 2.92 | 0.25 | 3.22 | 0.42 |
| PE (F)-claude-sonnet-4-0 | 0.25 | 0.43 | 3.04 | 0.43 | 3.57 | 0.43 |
| PE (F)-gemini-2.5-flash | 0.28 | 0.41 | 2.92 | 0.62 | 3.65 | 0.42 |
| PE (F+MBTI)-3.5-turbo | 0.25 | 0.40 | 2.97 | 0.35 | 2.55 | 0.40 |
| PE (F+MBTI)-4.1-mini | 0.24 | 0.44 | 3.20 | 0.37 | 3.32 | 0.41 |
| PE (F+MBTI)-4o-mini | 0.24 | 0.44 | 3.21 | 0.35 | 3.18 | 0.42 |
| PE (F+MBTI)-5-mini | 0.26 | 0.39 | 3.10 | 0.22 | 3.02 | 0.43 |
| PE (F+MBTI)-claude-sonnet-4-0 | 0.24 | 0.41 | 3.08 | 0.34 | 3.26 | 0.42 |
| PE (F+MBTI)-gemini-2.5-flash | 0.27 | 0.43 | 2.94 | 0.64 | 3.70 | 0.42 |

**(a) English**

**Performance Heatmap (Chinese)**

| Method-Base Model | M1 | M2 | M3 | M4 | M5 | M6 |
|---|---|---|---|---|---|---|
| LoRA-L3-ELYZA-8B | 0.94 | 0.64 | 2.58 | 0.56 | 3.54 | 0.79 |
| PE (0)-3.5-turbo | 1.04 | 0.42 | 3.05 | 0.25 | 2.19 | 0.46 |
| PE (0)-4.1-mini | 0.94 | 0.52 | 3.16 | 0.22 | 2.58 | 0.71 |
| PE (0)-4o-mini | 0.93 | 0.54 | 3.18 | 0.23 | 2.50 | 0.72 |
| PE (0)-5-mini | 1.04 | 0.38 | 3.20 | 0.05 | 2.05 | 0.60 |
| PE (0)-claude-sonnet-4-0 | 0.98 | 0.41 | 3.29 | 0.07 | 1.96 | 0.63 |
| PE (F)-3.5-turbo | 0.92 | 0.61 | 3.19 | 0.34 | 2.89 | 0.75 |
| PE (F)-4.1-mini | 0.93 | 0.61 | 3.21 | 0.30 | 3.12 | 0.76 |
| PE (F)-4o-mini | 0.91 | 0.61 | 3.19 | 0.33 | 3.16 | 0.74 |
| PE (F)-5-mini | 1.00 | 0.46 | 3.21 | 0.18 | 2.87 | 0.69 |
| PE (F)-claude-sonnet-4-0 | 0.89 | 0.63 | 3.22 | 0.44 | 3.75 | 0.79 |
| PE (F)-gemini-2.5-flash | 0.91 | 0.63 | 3.10 | 0.54 | 3.69 | 0.79 |
| PE (F+MBTI)-3.5-turbo | 0.95 | 0.56 | 3.08 | 0.30 | 2.60 | 0.70 |
| PE (F+MBTI)-4.1-mini | 0.92 | 0.62 | 3.20 | 0.33 | 3.23 | 0.77 |
| PE (F+MBTI)-4o-mini | 0.92 | 0.61 | 3.20 | 0.29 | 3.01 | 0.75 |
| PE (F+MBTI)-5-mini | 0.99 | 0.47 | 3.24 | 0.15 | 2.88 | 0.68 |
| PE (F+MBTI)-claude-sonnet-4-0 | 0.91 | 0.61 | 3.14 | 0.36 | 3.67 | 0.77 |
| PE (F+MBTI)-gemini-2.5-flash | 0.90 | 0.64 | 3.05 | 0.54 | 3.65 | 0.79 |

**(b) Chinese**

**Performance Heatmap (Japanese)**

| Method-Base Model | M1 | M2 | M3 | M4 | M5 | M6 |
|---|---|---|---|---|---|---|
| LoRA-L3-ELYZA-8B | 0.74 | 0.57 | 2.58 | 0.62 | 4.14 | 0.76 |
| PE (0)-3.5-turbo | 1.00 | 0.18 | 3.10 | 0.25 | 1.98 | 0.29 |
| PE (0)-4.1-mini | 0.72 | 0.31 | 3.32 | 0.25 | 2.90 | 0.59 |
| PE (0)-4o-mini | 0.72 | 0.28 | 3.25 | 0.24 | 2.57 | 0.58 |
| PE (0)-5-mini | 0.84 | 0.17 | 3.33 | 0.06 | 2.62 | 0.44 |
| PE (0)-claude-sonnet-4-0 | 0.81 | 0.18 | 3.21 | 0.13 | 2.02 | 0.47 |
| PE (0)-gemini-2.5-flash | 0.75 | 0.19 | 3.20 | 0.09 | 1.89 | 0.48 |
| PE (F)-3.5-turbo | 0.72 | 0.45 | 3.17 | 0.37 | 3.01 | 0.62 |
| PE (F)-4.1-mini | 0.70 | 0.42 | 3.29 | 0.31 | 3.99 | 0.64 |
| PE (F)-4o-mini | 0.69 | 0.46 | 3.20 | 0.35 | 4.22 | 0.66 |
| PE (F)-5-mini | 0.79 | 0.21 | 3.31 | 0.14 | 3.71 | 0.51 |
| PE (F)-claude-sonnet-4-0 | 0.68 | 0.54 | 3.09 | 0.51 | 4.49 | 0.71 |
| PE (F)-gemini-2.5-flash | 0.69 | 0.56 | 3.13 | 0.58 | 4.32 | 0.72 |
| PE (F+MBTI)-3.5-turbo | 0.73 | 0.44 | 3.19 | 0.37 | 3.27 | 0.62 |
| PE (F+MBTI)-4.1-mini | 0.70 | 0.45 | 3.23 | 0.33 | 4.20 | 0.65 |
| PE (F+MBTI)-4o-mini | 0.69 | 0.43 | 3.21 | 0.32 | 4.02 | 0.64 |
| PE (F+MBTI)-5-mini | 0.79 | 0.21 | 3.27 | 0.14 | 3.79 | 0.51 |
| PE (F+MBTI)-claude-sonnet-4-0 | 0.69 | 0.43 | 3.10 | 0.43 | 4.38 | 0.68 |
| PE (F+MBTI)-gemini-2.5-flash | 0.69 | 0.53 | 3.15 | 0.51 | 4.23 | 0.71 |

**(c) Japanese**

**Figure 3: Per-language performance heatmaps (compact layout). Rows are methods with base models; columns are metrics M1–M6 (M1 lower is better; others higher is better). Abbrev.: PE (0)=zero-shot, PE (F)=few-shot, PE (F+MBTI)=few-shot+MBTI, LoRA=PEFT (LoRA).**

**English.** Few-shot and MBTI-augmented prompting (**PE (F)**, **PE (F+MBTI)**) brighten M2–M3 columns, indicating stronger semantic and content alignment than zero-shot. `gpt-4.1-mini` and `gemini-2.5-flash` reach ~3.2–3.7 on M3/M5, showing balanced control of meaning and tone. LoRA maintains steady style imitation (M5=3.49, M6=0.45), suggesting robust consistency even without demonstrations.

**Chinese.** Performance varies more widely, with LoRA and few-shot MBTI variants leading stylistic metrics. `L3-ELYZA-8B` and `PE (F)-claude-sonnet-4-0` achieve top M5/M6 (≈3.7–3.8, 0.79), while **PE (F+MBTI)** balances substance and style (M2=0.62, M3=3.20, M6=0.77). High M1 (~0.9–1.0) suggests Chinese embeddings capture semantics less stably, making M2–M3 better indicators of alignment.

**Japanese.** Few-shot and MBTI methods dominate stylistic metrics (M5=4.0–4.5, M6>0.70), while M1–M3 remain moderate (~0.7, 3.2–3.3). This reflects Japanese morphology: looser lexical match but richer style control. LoRA and few-shot yield natural, coherent outputs once demonstrations are given.

**Cross-lingual Trends.** Few-shot prompting consistently enhances M2–M3 (semantic fidelity), MBTI inference improves M1/M3/M5 (personalized tone), and LoRA excels on M5–M6 (style/history). English models are stable across metrics; Chinese and Japanese show wider variance and greater headroom for personalization. Together, the heatmaps confirm that combining few-shot context and persona cues jointly optimizes both substance and style in multilingual personalization.

*5.2.3 Win–Tie–Lose Analysis Across Metrics.*

*Overall trends.* Figure 4 summarizes pairwise comparisons between methods for each metric. Across the board, **few-shot prompting outperforms zero-shot prompting**: PE (0) is dominated by "Lose" shares on most metrics (e.g., ~73–80% on M2/M4/M6), indicating that zero-shot prompting alone is insufficient for stable personalization.

*Metric-specific observations.*

- **Semantic and structure-oriented metrics favor few-shot.** On **M2** (sentence-level similarity), PE (F) attains the highest win share (about **46.7%**). PE (F) is also strongest on

**M4** (length/format consistency) and **M6** (similarity to user history), with wins around **53.3%**.

- **Personality cues add value for style/content choices.** Adding MBTI inference, **PE (F+MBTI)**, yields the top win share on **M1** (distance; **40.0%**, lower is better), **M3** (content similarity by an LLM; **33.3%**), and **M5** (style similarity; **40.0%**). This suggests that coarse personality hypotheses help models better match users' lexical and stylistic preferences.

- **LoRA is competitive but not dominant.** LoRA achieves wins of roughly ~20% on several metrics, indicating that parameter-efficient tuning is effective in specific settings (e.g., more stable formatting or style) but is not the overall leader in cross-language, multi-metric comparisons.
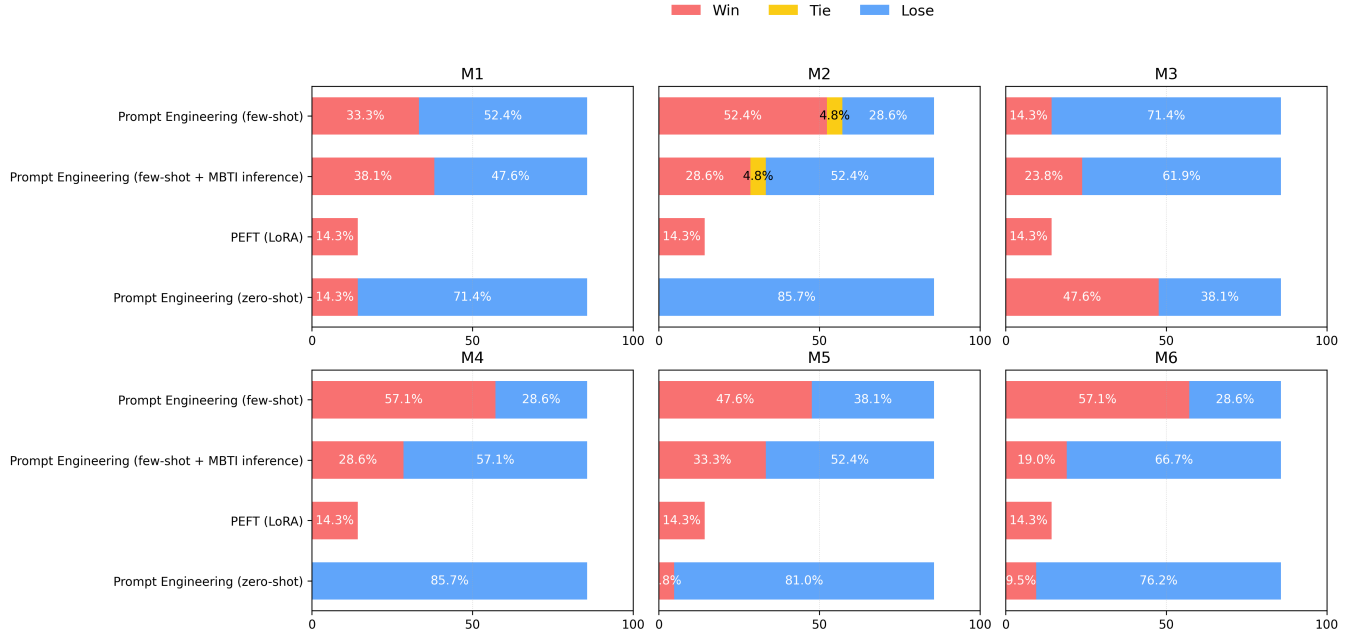
*Practical takeaways.* For *semantic fidelity* (M2, M3) and *alignment to personal history* (M6), prefer **few-shot prompting**, and include **MBTI inference** when available to further improve content and style matching (especially on M1/M3/M5). When *format control* (M4) is critical, **PE (F)** is the most stable choice. No single method dominates every metric; the selection should be *metric-driven* and conditioned on the available user evidence (demonstrations vs. lightweight tuning).

## 6 CONCLUSION

We presented **Your Next Token Prediction**, a new benchmark and framework for evaluating personalized response generation in large language models. By integrating multi-day human–agent interactions with MBTI-grounded personality modeling, our system enables the collection of fine-grained user–specific dialogue data. The released 100-user multilingual dataset and benchmark provide a foundation for studying the semantic and stylistic alignment of LLMs to individual users. Through extensive evaluations of prompt engineering and fine-tuning methods, we demonstrate both the challenges and potential of achieving truly personalized alignment. Future work will explore scalable adaptation strategies that combine efficiency, privacy, and long-term user modeling.

**Figure 4: Win/Tie/Lose across metrics (M1–M6) over all languages and base models. Each panel summarizes the percentage of configurations a method *wins/ties/loses* on that metric. Few-shot methods dominate zero-shot; adding MBTI nudges wins on M1/M3/M5; LoRA is competitive on style (M5) and history (M6).**

## ACKNOWLEDGMENTS

If you wish to include any acknowledgments in your paper (e.g., to people or funding agencies), please do so using the 'acks' environment. Note that the text of your acknowledgments will be omitted if you compile your document with the 'anonymous' option.

## REFERENCES

[1] Louis Castricato, Nathan Lile, Rafael Rafailov, Jan-Philipp Fränken, and Chelsea Finn. 2024. Persona: A reproducible testbed for pluralistic alignment. *arXiv preprint arXiv:2407.17387* (2024).

[2] Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai, and Zuozhu Liu. 2024. Pad: Personalized alignment of llms at decoding-time. *arXiv preprint arXiv:2410.04070* (2024).

[3] Konstantina Christakopoulou, Alberto Lalama, Cj Adams, Iris Qu, Yifat Amir, Samer Chucri, Pierce Vollucci, Fabio Soldo, Dina Bseiso, Sarah Scodel, et al. 2023. Large language models for user interest journeys. *arXiv preprint arXiv:2305.15498* (2023).

[4] Personalized Parameter-Efficient Fine-tuning. [n.d.]. Democratizing Large Language Models via Personalized Parameter-Efficient Fine-tuning. ([n. d.]).

[5] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR* 1, 2 (2022), 3.

[6] Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *arXiv preprint arXiv:2310.11564* (2023).

[7] Hannah Rose Kirk, Alexander Whitefield, Paul Rottger, Andrew M Bean, Katerina Margatina, Rafael Mosquera-Gomez, Juan Ciro, Max Bartolo, Adina Williams, He He, et al. 2024. The PRISM alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models. *Advances in Neural Information Processing Systems* 37 (2024), 105236–105344.

[8] Ishita Kumar, Snigdha Viswanathan, Sushrita Yerra, Alireza Salemi, Ryan A Rossi, Franck Dernoncourt, Hanieh Deilamsalehy, Xiang Chen, Ruiyi Zhang, Shubham Agarwal, et al. 2024. Longlamp: A benchmark for personalized long-form text generation. *arXiv preprint arXiv:2407.11016* (2024).

[9] Xinyu Li, Ruiyang Zhou, Zachary C Lipton, and Liu Leqi. 2024. Personalized language modeling from personalized human feedback. *arXiv preprint arXiv:2402.05133* (2024).

[10] Jiongnan Liu, Yutao Zhu, Shuting Wang, Xiaochi Wei, Erxue Min, Yu Lu, Shuaiqiang Wang, Dawei Yin, and Zhicheng Dou. 2024. Llms+ persona-plug= personalized llms. *arXiv preprint arXiv:2409.11901* (2024).

[11] Thao Nguyen, Krishna Kumar Singh, Jing Shi, Trung Bui, Yong Jae Lee, and Yuheng Li. 2025. Yo'Chameleon: Personalized Vision and Language Generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*. 14438–14448.

[12] Chris Richardson, Yao Zhang, Kellen Gillespie, Sudipta Kar, Arshdeep Singh, Zeynab Raeesy, Omar Zia Khan, and Abhinav Sethy. 2023. Integrating summarization and retrieval for enhanced personalization via large language models. *arXiv preprint arXiv:2310.20081* (2023).

[13] Alireza Salemi, Surya Kallumadi, and Hamed Zamani. 2024. Optimization methods for personalizing large language models through retrieval augmentation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 752–762.

[14] Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2023. Lamp: When large language models meet personalization. *arXiv preprint arXiv:2304.11406* (2023).

[15] Krishna Sayana, Raghavendra Vasudeva, Yuri Vasilevski, Kun Su, Liam Hebert, James Pine, Hubert Pham, Ambarish Jash, and Sukhdeep Sodhi. 2025. Beyond Retrieval: Generating Narratives in Conversational Recommender Systems. In *Companion Proceedings of the ACM on Web Conference 2025*. 2411–2420.

[16] Xiangru Tang, Xingyao Zhang, Yanjun Shao, Jie Wu, Yilun Zhao, Arman Cohan, Ming Gong, Dongmei Zhang, and Mark Gerstein. 2024. Step-back profiling: Distilling user history for personalized scientific writing. *arXiv preprint arXiv:2406.14275* (2024).

[17] Shujin Wu, May Fung, Cheng Qian, Jeonghwan Kim, Dilek Hakkani-Tur, and Heng Ji. 2024. Aligning llms with individual preferences via interaction. *arXiv preprint arXiv:2410.03642* (2024).

[18] Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. 2023. Flask: Fine-grained language model evaluation based on alignment skill sets. *arXiv preprint arXiv:2307.10928* (2023).

[19] Siyan Zhao, Mingyi Hong, Yang Liu, Devamanyu Hazarika, and Kaixiang Lin. 2025. Do llms recognize your preferences? evaluating personalized preference following in llms. *arXiv preprint arXiv:2502.09597* (2025).

[20] Thomas P Zollo, Andrew Wei Tung Siah, Naimeng Ye, Ang Li, and Hongseok Namkoong. 2024. Personalllm: Tailoring llms to individual preferences. *arXiv preprint arXiv:2409.20296* (2024).