

# Persona-Aware Alignment Framework for Personalized Dialogue Generation

Guanrong Li Xinyu Liu Zhen Wu\* Xinyu Dai

National Key Laboratory for Novel Software Technology, Nanjing University  
{grli, xinyuliu}@smail.nju.edu.cn, {wuz, daixinyu}@nju.edu.cn

## Abstract

Personalized dialogue generation aims to leverage persona profiles and dialogue history to generate persona-relevant and consistent responses. Mainstream models typically rely on token-level language model training with persona dialogue data, such as Next Token Prediction, to implicitly achieve personalization, making these methods tend to neglect the given personas and generate generic responses. To address this issue, we propose a novel Persona-Aware Alignment Framework (PAL), which directly treats persona alignment as the training objective of dialogue generation. Specifically, PAL employs a two-stage training method including *Persona-aware Learning* and *Persona Alignment*, equipped with an easy-to-use inference strategy *Select then Generate*, to improve persona sensitivity and generate more persona-relevant responses at the semantics level. Through extensive experiments, we demonstrate that our framework outperforms many state-of-the-art personalized dialogue methods and large language models.

## 1 Introduction

Building personalized dialogue agents capable of human-like interactions has emerged as a critical research frontier with profound societal implications. Unlike task-oriented chatbots that only prioritize functional efficiency, agents imbued with consistent personas enable applications that demand empathy, rapport, and contextual adaptability. For instance, in mental health support systems, chatbots adopting a compassionate persona can foster user trust for more effective emotional counseling (Sarikaya et al., 2016). In education, tutors with tailored personalities improve student engagement through relatable interactions (McTear, 2022). Personalized agents also unlock scalable

solutions for senior companionship, customer service avatars, and interactive storytelling, where alignment with user expectations hinges on maintaining coherent personas (Singh and Beniwal, 2022). Therefore, how to build personalized dialogue agents has drawn increasing attention in recent years (Sarikaya et al., 2016; McTear, 2022; Singh and Beniwal, 2022). Many studies integrate predefined descriptions of specific individuals, often referred to as personas, with dialogue history to generate personalized responses (Pei et al., 2021; Wu et al., 2021; Tang et al., 2023; Han et al., 2024). As shown in Figure 1, given a dialogue history and a personas profile, the upper response “I do not have time to watch TV, I am an attorney so i work a lot.” aligns with personas better than the bottom response. Here, alignment refers to the generated responses reflecting predefined persona attributes.

Mainstream works usually use persona and dialogue history as the input of models and strongly depend on data driving to generate personalized responses (Song et al., 2021; Liu et al., 2022; Chen et al., 2023a). Although achieving remarkable progress, these methods often struggle to generate responses that align with the given personas. A primary reason for this limitation lies in the token-level training objectives of these models, such as Next Token Prediction (NTP) and Mask Prediction (MP). These objectives are primarily designed for language modeling, which excels at producing grammatically correct and contextually coherent sentences, while falling short in obtaining good persona alignment.

To address this issue, we propose learning alignment explicitly between generated responses and given personas as the optimization objective. Unlike token-level language modeling objectives, this response-level alignment objective captures personas and generates personalized responses at the semantic level. Two practical challenges re-

\* Corresponding author.

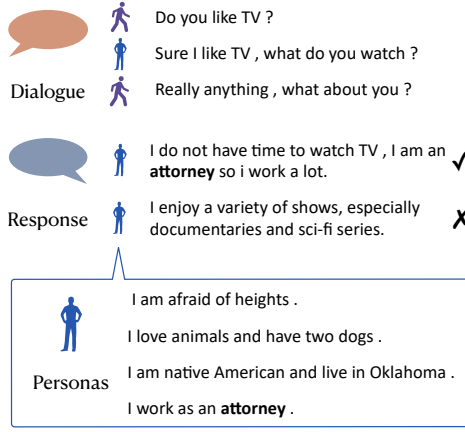


Figure 1: An example for the personalized dialogue. ✓ represents the persona consistent response and ✗ represents a generic response instead of a personal response.

main to achieve this goal. First, quantifying the degree of alignment of personas as a learning signal is challenging because persona alignment requires strong semantic understanding and lacks well-established calculation metrics. Second, a persona profile contains multiple irrelevant persona descriptions, typically only one is relevant to the current dialogue history. The irrelevant personas could distract or mislead the models.

To overcome the challenges, we introduce a novel **Persona-Aware Alignment Framework (PAL)**, which treats persona alignment as the primary optimization objective. PAL employs a two-stage training approach and a *Select then Generate* inference strategy to improve persona sensitivity and generate persona-relevant responses. In the first training stage, called *Persona-aware Learning*, we address the challenge of irrelevant personas by designing a mixed training task including two subtasks: *Dialogue-Informed Persona Selection* and *Persona-Enhanced Dialogue Generation*. They aim respectively to learn which persona is relevant to the current dialogue and generate the persona-aware response. In the second training stage, inspired by preference learning (Xiong et al., 2024), we propose *Persona Alignment* with Direct Preference Optimization (DPO) (Rafailov et al., 2023) to address the quantification and explicit learning of persona alignment. As there is no available paired data, we construct the training data by pairing contrasting generated responses with ground-truth responses. Finally, we naturally employ a *Select then Generate* inference strategy to filter irrelevant personas. Similar to the per-

sona extractor proposed in Xu et al. (2022), our approach first selects the most contextually relevant personas and then generates a response conditioned on the selected information.

Our contributions can be summarized as:

- We propose the learning of persona alignment as the optimization objective for personalized dialogue. To our knowledge, it is the first to explicitly learn persona information from a response level objective.
- We introduce the novel **Persona-Aware Alignment Framework (PAL)** to tackle specific challenges of persona alignment learning. This framework helps improve persona sensitivity of models and generate persona-relevant responses.
- We conduct comprehensive experiments across different foundation models, datasets, and languages. The superior performance of our framework over both state-of-the-art baselines and well-known large language models demonstrates its effectiveness and broad applicability.

## 2 Persona-Aware Alignment Framework

In this section, we will introduce our **Persona-Aware Alignment Framework (PAL)** in detail. Formally, for a user  $u \in \mathcal{U}$ , there is a set of personas denoted by  $P = \{p_1, p_2, \dots, p_l\}$  and an associated set of multi-turn dialogue contexts resulting from interactions with other users, represented as  $C = \{(q_1, r_1), (q_2, r_2), \dots, (q_n, r_n)\}$ . In these contexts,  $q_i$  denotes the dialogue utterance from another user, and  $r_i$  is the corresponding response from  $u$ , where  $n$  is the total number of dialogue turns. Given the next conversation turn  $q_{n+1}$ , the goal of the personalized dialogue agent is to generate a personalized response  $r_{n+1}$  that aligns with user  $u$ 's personas. PAL utilizes semantic-level persona alignment as the optimization target, aiming at training the personalized dialogue agent to generate responses aligned with given personas. Figure 2 shows the structure of PAL. It consists of a two-stage training method, i.e., *Persona-aware Learning* and *Persona Alignment*, and a *Select then Generate* inference strategy.

### 2.1 Persona-aware Learning

As previously mentioned in the Introduction, the persona profile usually has one most relevant per-

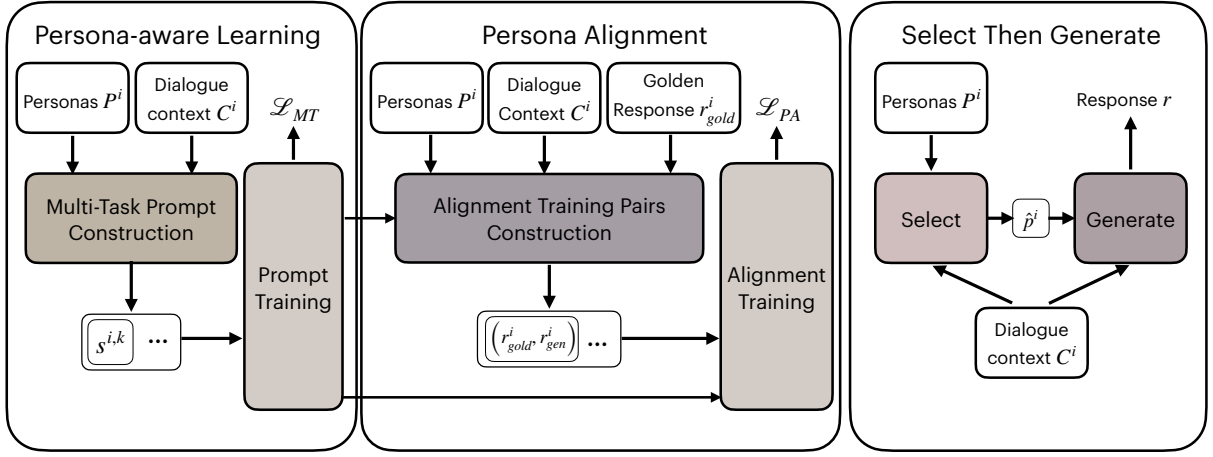


Figure 2: The overview of our Persona-Aware Alignment Framework (PAL) includes a two-stage training strategy: (1) *Persona-aware Learning*, and (2) *Persona Alignment*, as well as a *Select then Generate* inference strategy. The arrows trace the flow of information, showing how each stage converts its inputs into outputs. In the *Persona-aware Learning* stage, the inputs are the persona descriptions  $P^i$  and dialogue context  $C^i$ . A multi-task prompt-construction module turns these inputs into training samples  $s^{i,k}$ , which are used for prompt tuning, with the training loss  $\mathcal{L}_{MT}$  as the output of this stage. In the *Persona Alignment* stage, the inputs include persona descriptions  $P^i$ , dialogue context  $C^i$ , and the gold response  $r_{gold}^i$ . An alignment-pair constructor forms  $(r_{gold}^i, r_{gen}^i)$ , where  $r_{gen}^i$  is produced by the model from the previous stage. These pairs yield the alignment loss  $\mathcal{L}_{PA}$ , the output of this stage by Alignment Training. In the *Select-then-Generate* inference strategy, the inputs are the persona descriptions  $P^i$  and dialogue context  $C^i$ . A selection module picks the persona  $\hat{p}^i$  most relevant to the context. A response generator then produces the final reply  $r$ , explicitly highlighting the selected persona.

sona description for the dialogue history, and the others could be noises misleading the models. To bridge this gap, we propose a mixed training method *Persona-aware Learning*. It involves two subtasks, *Dialogue-Informed Persona Selection* and *Persona-Enhanced Dialogue Generation*. *Dialogue-Informed Persona Selection* aims to teach the model which persona is relevant through explicit persona learning, which improves the sensitivity of the model to personas. *Persona-Enhanced Dialogue Generation* aims to generate responses based on given personas and dialogue history. The two subtasks have different input and output formats. To reduce semantic shift from task formats during mixed training, we construct prompts for each task and transform them into a unified natural language format.

### 2.1.1 Task Design

Given a user  $u$  with personas  $P = \{p_1, p_2, \dots, p_l\}$  and a dialogue history  $C$ , the original formulations of the two key subtasks are as follows.

**Dialogue-Informed Persona Selection.** This subtask is designed to enhance the model’s persona understanding ability by selecting the most relevant persona based on the dialogue history. The

goal is to identify the persona that best aligns with the dialogue history  $C$  from personas  $P$ , which is typically formulated as a classification problem:

$$\hat{p} = \arg \max_{p \in P} \text{sim}(f(C), g(p)), \quad (1)$$

where  $\hat{p}$  denotes the selected persona that is most relevant to the dialogue context  $C$ . Here,  $f(C), g(p)$  are encoding functions that map the dialogue history and persona into semantic representations, respectively. The function  $\text{sim}(\cdot)$  computes the similarity between these representations.

**Persona-Enhanced Dialogue Generation.** We design this subtask to enhance the model’s personalized generation ability. It focuses on generating an appropriate response for the dialogue. It is typically formulated as a token prediction task:

$$\hat{r} = \arg \max_r \prod_{t=1}^T \mathcal{P}(w_t | w_{<t}, C, P) \quad (2)$$

where  $\hat{r}$  represents the generated response that best fits the dialogue context,  $r = \{w_1, w_2, \dots, w_T\}$  is the sequence of tokens in the response, with  $T$  denoting the length of the response,  $w_t$  is the token at position  $t$  in the response, and  $w_{<t} = \{w_1, \dots, w_{t-1}\}$  represents the

sequence of all preceding tokens, and  $\mathcal{P}(\cdot)$  denotes the probability of predicting a token.

### 2.1.2 Prompt Construction and Training

**Multi-Task Prompt Construction.** As mentioned above, these subtasks have different input and output formats, which require specialized modules and limit generality. Therefore, we transform them into unified natural language formats through constructing prompts for each task, and integrate them into one unified task with supervised fine-tuning. The corresponding prompts are shown as follows:

**Dialogue-Informed Persona Selection.**

The user’s persona is described with: <personas>.

If a persona description is required to generate a response, select the most appropriate one. If no persona is needed, respond with ‘No persona data needed’.

Dialogue context: <dialogue context>.

The preferred persona is: <related persona>.

**Persona-Enhanced Dialogue Generation.**

The user’s persona is described with: <personas>.

Please generate a response to the dialogue.

Dialogue context: <dialogue context>.

Response: <response>

where “<personas>” in the table represents the personas  $P$ , the “<dialogue context>” represents the corresponding dialogue history  $C$ , the “<related persona>” represents the selected persona  $\hat{p}$  in Equation 1, and “<response>” represents the generated response  $\hat{r}$  in Equation 2.

**Prompt Training.** After constructing the prompts, we employ next-token prediction (NTP) to facilitate multi-task training. The combined loss function for mixed training is defined as:

$$\mathcal{L}_{\text{Mix}}(\theta) = - \sum_k \sum_{t=1}^{T^k} \log \mathcal{P}(s_t^k | m_k, s_{<t}^k; \theta)$$

where  $\theta$  represents the model parameters,  $k$  represents the index over tasks, with  $k \in \{1, 2\}$ , corresponding to Task 1: *Dialogue-Informed Persona Selection* and Task 2: *Persona-Enhanced Dialogue Generation*.  $m_k$  represents the prompt constructed for task  $k$ .  $s^k$  is the target sequence in task  $k$ , and  $t$  is the time step index within the

target sequence for each task.  $T^k$  represents the total number of tokens in the target sequence  $s^k$ .

## 2.2 Persona Alignment

One innovation of this work is to train the personalized dialogue agent by directly aligning the generated responses with the given personas. The key challenge in this process is quantifying this alignment because there is currently no meaningful metric available. Given the strong fitting capabilities of neural networks, we leverage neural networks to estimate the alignment between responses and personas. Specifically, we first construct persona alignment training pairs. Then, we can train a Proximal Policy Optimization (PPO) (Schulman et al., 2017) pipeline to obtain an alignment evaluation model. In practice, we instead adopt the Direct Preference Optimization (DPO) method (Rafailov et al., 2023). DPO allows us to implicitly incorporate the alignment measurement model and avoids the overhead associated with reinforcement learning methods like PPO, offering a simpler, more efficient approach.

### 2.2.1 Alignment Training Pairs Construction

Since there are no ground-truth alignment scores available, we construct a learning-to-rank task inspired by Preference Learning methods (Xiong et al., 2024). Our training objective is to assign higher scores to responses that better align with the given personas. Specifically, we regard the ground truth response  $r_{\text{gold}}$  from personalized dialogue datasets as the chosen response of the Direct Preference Optimization (DPO) (Rafailov et al., 2023). Considering that there is no manually labeled rejected response, we adopt an alternative approach and utilize the trained model from the *Persona-aware Learning* stage to generate responses  $r_{\text{gen}}$  based solely on the dialogue history, excluding personal information, as rejected responses. The prompts are shown as follows:

**Prompt for Alignment Training Pairs Construction**

Please generate a response to the dialogue.

Dialogue context: <dialogue context>.

Response: <response>

In general, we believe that the golden responses  $r_{\text{gold}}$  provided by the personalized dialogue datasets are more aligned with the personas than the generated responses  $r_{\text{gen}}$  without persona



input. Therefore, we construct data pairs  $\mathcal{D} = \{(r_{\text{gold}}, r_{\text{gen}})\}$  for alignment training by pairing the golden responses with the less-aligned generated responses.

### 2.2.2 Alignment Training

This stage aims to directly tune the model guided by persona alignment. As mentioned above, we utilize the DPO method for alignment training. The training targets are defined as follows:

$$\delta_{\text{gold}} = \log \frac{\pi_{\theta}(r_{\text{gold}} | C, P)}{\pi_{\theta_{\text{ref}}}(r_{\text{gold}} | C, P)}$$

$$\delta_{\text{gen}} = \log \frac{\pi_{\theta}(r_{\text{gen}} | C, P)}{\pi_{\theta_{\text{ref}}}(r_{\text{gen}} | C, P)}$$

$$\mathcal{L}_{PA} = -\mathbb{E}_{(r_{\text{gold}}, r_{\text{gen}}) \sim \mathcal{D}} [\log \sigma(\beta(\delta_{\text{gold}} - \delta_{\text{gen}}))]$$

where  $\pi_{\theta}$  represents the model being trained, parameterized by  $\theta$ .  $\pi_{\theta_{\text{ref}}}$  is the reference model with parameters  $\theta_{\text{ref}}$ , which remains fixed during this stage.  $\mathcal{D}$  represents the set of training pairs.  $\beta$  is a scaling factor that adjusts the sharpness of the preference,  $\sigma$  denotes the sigmoid function.

### 2.3 Select then Generate Strategy

To address the challenge of persona relevance and minimize the influence of noisy personas at the inference stage, we propose a *Select then Generate* strategy. This approach begins by selecting the most relevant persona based on the dialogue history using a carefully designed prompt. By focusing on the most pertinent persona, we reduce the risk of irrelevant or misleading information affecting the response generation. After selecting the most relevant persona, we construct a new prompt that highlights this persona, guiding the response generation process. This ensures that the generated responses are both contextually appropriate and closely aligned with the user’s specific persona attributes. The prompts are shown as follows:

#### Prompt for Select

The user’s persona is described with: <personas>.  
If a persona description is required to generate a response, select the most appropriate one. If no persona is needed, respond with ‘No persona data needed’.  
Dialogue context: <dialogue context>.  
The preferred persona is: <related persona>.

|          | PERSONA-CHAT | Baidu-Persona-Chat |
|----------|--------------|--------------------|
| Language | English      | Chinese            |
| Train    | 65719        | 353016             |
| Valid#1  | 6500         | 14000              |
| Valid#2  | 1301         | 3027               |
| Test     | 7512         | 5929               |

Table 1: Number of samples in each dataset.

#### Prompt for Generate

The user’s persona is described with: <personas>.  
The most related persona is <related persona>.  
Please generate a response to the dialogue.  
Dialogue context: <dialogue context>.  
Response: <response>

## 3 Experimental Settings

### 3.1 Datasets

To assess the effectiveness and adaptability of our framework, we conducted comprehensive experiments using two widely recognized datasets across different languages: PERSONA-CHAT (Zhang et al., 2018a) and Baidu-Persona-Chat<sup>1</sup>. PERSONA-CHAT is an English dataset comprising an extensive collection of dialogues paired with corresponding personas. We conducted experiments on both the original and revised versions of this dataset. In the revised version, the personal information has been rewritten. Both versions share the same test set. Baidu-Persona-Chat is a Chinese personalized dialogue dataset similar in format to PERSONA-CHAT. We split the datasets into four parts (Train, Valid#1, Valid#2, and Test) to meet the specific requirements of the personalized dialogue task. The data statistics are shown in Table 1. During the *Persona-aware Learning* stage, we use both validation subsets (Valid #1 and Valid #2) for evaluation. In the *Persona Alignment* stage, Valid #1 is utilized for training and Valid #2 for validation. This partitioning approach prevents data leakage and is supported by prior research, as well as our experimental findings, which underscore the effectiveness of using relatively smaller datasets for alignment tasks (Ouyang et al., 2022; Lee et al., 2023).

<sup>1</sup><https://aistudio.baidu.com/datasetdetail/351937>

### 3.2 Evaluation Metrics

We evaluate the model using both automatic and human metrics. For the automatic metrics, we adopt **BLEU** (Papineni et al., 2002) and **ROUGE** (Lin, 2004), both of which measure the token overlap. Additionally, we utilize **Entropy** (Zhang et al., 2018b), which measures the diversity of personalized responses. Following previous works (Madotto et al., 2019), we also adopt the **C.score** to assess the consistency between the model-generated responses and the provided personas. The methodology for calculating the C.score is detailed as follows:

$$NLI(p_l, r_n) = \begin{cases} 1 & \text{if } p_l \text{ entails } r_n \\ 0 & \text{if } p_l \text{ is independent of } r_n \\ -1 & \text{if } p_l \text{ contradicts } r_n \end{cases}$$

$$\text{C.score} = \sum_{p_l^i \in P^i} NLI(p_l^i, r_n^i)$$

where  $P = \{p_1, p_2, \dots, p_l\}$  represents the personas of user  $u$ ,  $r_n$  denotes the generated response of the last dialogue turn.  $l$  indexes the personas and  $n$  indicates the dialogue turn. We utilize a Natural Language Inference (NLI) model to calculate the Consistency score (C.score). Specifically, we utilized a RoBERTa (Liu et al., 2019) model designed for three-way classification tasks (entailment, neutral, and contradiction). For the English experiments, we further trained this RoBERTa model on SNLI (Glockner et al., 2018) and MultiNLI (Kim et al., 2019) datasets. For the Chinese experiments, we trained this RoBERTa model on the OCNLI (Hu et al., 2020) and CMNLI (Xu et al., 2020) datasets. We then further trained the NLI model on the PERSONA-CHAT and Baidu-Persona-Chat datasets. The NLI model achieve accuracies of 84.1% and 82.3%, respectively.

For human evaluation, we randomly select 100 samples from the test set. Two annotators were asked to evaluate responses based on three aspects: (1) Fluency, (2) Coherence, and (3) Persona Consistency. The assigned scores of 1, 2, 3 correspond to unacceptable, acceptable, and satisfactory levels, respectively. The inter-annotator agreement for our study was measured by the Fleiss Kappa coefficient, with a score of 0.62, indicating substantial agreement. More details about the annotations are shown in Appendix A.

### 3.3 Baselines

We apply our PAL framework to three foundation models to evaluate its effectiveness.

- **GPT-2** (Radford et al., 2019): A pre-trained language model developed by OpenAI, built on the Transformer architecture. We utilize the base version of this model series.
- **DialoGPT** (Zhang et al., 2020): A variant of GPT-2 specifically adapt for dialogue generation, developed by Microsoft. We utilize the base version. As there is no Chinese-pretrained version of DialoGPT, we only conduct experiments on the PERSONA-CHAT dataset.
- **Llama 3.1 8B** (Dubey et al., 2024): An open-source large language model developed by Meta, known for its strong capabilities and wide applicability. We employ the 8B version to highlight the performance of our PAL framework on a large-scale model.

We compare several different strategies and some state-of-the-art (SOTA) methods.

- **Fine-Tuning** (Kenton and Toutanova, 2019): Fine-Tuning is a widely used strategy to adapt a pre-trained model to a specific task. In this work, we continue training the pre-trained model on personalized dialogue datasets. Due to the high computational cost and insufficient data for training large language models, we conducted experiments only on GPT-2 and DialoGPT.
- **Few-Shot Prompting** (Brown et al., 2020): In this approach, the model is provided with a small number of demonstrations during inference to guide its responses. We include two demonstration dialogues in our experiments that incorporate persona information in the input prompt, helping the model generate responses that are more aligned with the user’s persona. This strategy is utilized only with Llama 3.1 8B.
- **Prompt Tuning** (Lester et al., 2021): The Prompt Tuning technique adapts a pre-trained model to specific tasks by introducing a small set of task-specific parameters.
- **SimOAP** (Zhou et al., 2023b): It adopts a two-stage strategy for personalized dialogue that involves oversampling and post-evaluation during generation.

| Models         | Original PERSONA-CHAT |                 |                 |               |                 | Revised PERSONA-CHAT |                 |                 |               |                |
|----------------|-----------------------|-----------------|-----------------|---------------|-----------------|----------------------|-----------------|-----------------|---------------|----------------|
|                | BLEU-1                | BLEU-2          | ROUGE-L         | Entropy       | C.score         | BLEU-1               | BLEU-2          | ROUGE-L         | Entropy       | C.score        |
| ORIG           | <u>13.97</u>          | <b>7.40</b>     | <u>16.20</u>    | <u>6.55</u>   | <b>0.696</b>    | 12.20                | 6.40            | <b>15.92</b>    | <u>5.75</u>   | <u>0.313</u>   |
| CLV            | <b>16.79</b>          | <u>6.74</u>     | <u>15.02</u>    | <b>8.30</b>   | 0.576           | <b>16.65</b>         | <u>6.75</u>     | <u>15.05</u>    | <b>8.34</b>   | <b>0.416</b>   |
| LMEDR          | 12.56                 | 6.60            | <b>16.89</b>    | 5.42          | <u>0.582</u>    | <u>14.34</u>         | <b>6.91</b>     | 13.72           | 5.47          | 0.312          |
| GPT-3.5        | 2.12                  | 0.71            | 6.19            | <b>8.65</b> † | 0.809           | 2.12                 | 0.71            | 6.19            | <b>8.65</b> † | 0.809          |
| +Few-Shot      | <b>6.20</b>           | <u>0.79</u>     | <b>8.61</b>     | 5.01          | <b>0.884</b>    | <b>6.20</b>          | <u>0.79</u>     | <b>8.61</b>     | 5.01          | <b>0.884</b> † |
| GPT 4o mini    | 3.00                  | <b>0.99</b>     | <u>8.09</u>     | <u>7.62</u>   | 0.670           | 3.00                 | <b>0.99</b>     | <u>8.09</u>     | <u>7.62</u>   | 0.670          |
| +Few-Shot      | 2.22                  | 0.07            | 5.39            | 5.33          | 0.825           | 2.22                 | 0.07            | 5.39            | 5.33          | 0.825          |
| Gemini Flash   | 1.11                  | 0.36            | 4.95            | 6.95          | 0.641           | 1.11                 | 0.36            | 4.95            | 6.95          | 0.641          |
| +Few-Shot      | <u>3.31</u>           | 0.67            | 3.11            | 5.32          | <u>0.870</u>    | <u>3.31</u>          | 0.67            | 3.11            | 5.32          | <u>0.870</u>   |
| GPT-2          | 5.97                  | 2.50            | 9.74            | 5.18          | 0.333           | 5.97                 | 2.50            | 9.74            | 5.18          | 0.333          |
| +FineTuning    | 13.10                 | 7.06            | 15.92           | 6.08          | 0.173           | 12.73                | 6.62            | <u>15.34</u>    | 6.28          | 0.103          |
| +Prompting     | <u>13.49</u>          | <u>7.52</u>     | <u>16.35</u>    | 6.45          | 0.454           | <u>13.25</u>         | <u>7.09</u>     | <b>16.20</b>    | 6.26          | 0.218          |
| +SimOAP        | 8.09                  | 2.26            | 8.42            | <b>8.47</b> * | 0.170           | 9.00                 | 3.19            | 9.23            | <b>8.39</b>   | 0.153          |
| +SPT           | 5.55                  | 2.65            | 9.82            | 6.01          | <u>0.799</u>    | 4.84                 | 2.20            | 9.03            | 4.97          | <b>0.440</b>   |
| +PAL(ours)     | <b>17.05</b> *        | <b>8.77</b> *   | <b>16.66</b> *  | <u>6.93</u>   | <b>0.811</b> *  | <b>15.94</b> *       | <b>7.89</b> *   | 14.55           | <u>6.75</u>   | <u>0.427</u>   |
| DialoGPT       | 5.09                  | 2.00            | 6.16            | 6.59          | 0.485           | 5.09                 | 2.00            | 6.16            | 6.59          | <b>0.485</b>   |
| +FineTuning    | 12.89                 | 6.84            | <u>15.48</u>    | 6.10          | 0.125           | 12.48                | 6.57            | <u>15.65</u>    | 6.30          | 0.096          |
| +Prompting     | <u>13.60</u>          | <u>7.36</u>     | <b>16.42</b>    | 6.54          | 0.445           | <u>13.39</u>         | <b>7.20</b>     | <b>16.42</b>    | 6.36          | 0.209          |
| +SimOAP        | 9.17                  | 3.36            | 9.23            | <b>8.39</b>   | 0.085           | 9.19                 | 3.25            | 9.28            | <b>8.39</b>   | 0.080          |
| +SPT           | 4.83                  | 2.26            | 8.59            | 5.60          | <u>0.509</u>    | 4.78                 | 1.44            | 5.42            | 4.82          | 0.303          |
| +PAL(ours)     | <b>15.11</b> *        | <b>7.65</b> *   | 15.38           | <u>6.81</u>   | <b>0.585</b> *  | <b>13.69</b> *       | <u>7.13</u>     | 15.61           | <u>6.69</u>   | <u>0.349</u>   |
| Llama 3.1 - 8B | 3.11                  | 0.98            | 5.78            | 5.56          | 0.611           | 3.11                 | 0.98            | 5.78            | 5.56          | <u>0.611</u>   |
| +Few-Shot      | 4.89                  | 1.97            | 9.40            | 4.28          | 0.219           | 4.89                 | 1.97            | 9.40            | 4.28          | 0.219          |
| +Prompting     | <u>18.95</u>          | <u>8.27</u>     | <u>16.70</u>    | 5.77          | 0.642           | <u>18.24</u>         | <u>7.64</u>     | <u>16.09</u>    | 5.78          | 0.466          |
| +SimOAP        | 12.69                 | 4.11            | 11.44           | 6.74          | 0.680           | 12.69                | 3.96            | 11.43           | 5.63          | 0.441          |
| +SPT           | 8.61                  | 3.81            | 8.14            | <u>6.98</u>   | <u>0.694</u>    | 7.61                 | 2.71            | 7.80            | <u>6.27</u>   | 0.594          |
| +PAL(ours)     | <b>25.12</b> *†       | <b>14.00</b> *† | <b>23.88</b> *† | <b>7.27</b> * | <b>0.909</b> *† | <b>24.89</b> *†      | <b>13.42</b> *† | <b>22.75</b> *† | <b>7.33</b> * | <b>0.625</b> * |

Table 2: Automatic Evaluation Results on the PERSONA-CHAT Dataset. The best results are indicated in **bold**, while the secondary results are marked with underlined. \* denotes statistically significant ( $p < 0.05$ ) improvements over baselines on the same foundation models. † indicates superior performance over all other results ( $p < 0.05$ ). All significance testing was conducted using independent samples t-tests.

- **SPT** (Huang et al., 2024): SPT adapts a pre-trained model for personalized conversations by introducing soft prompts. Instead of fine-tuning the entire model, it uses a trainable retriever to selectively choose suitable soft prompts based on the input context.

Additionally, we compare with SOTA methods that have specifically designed modules.

- **ORIG** (Chen et al., 2023a): ORIG addresses the issue of persona order sensitivity in personalized dialogue generation.
- **CLV** (Tang et al., 2023): This method utilizes both sparse and dense representations of personas for dialogue generation.
- **LMEDR** (Chen et al., 2023b): It learns to memorize entailment and discourse relations for persona-consistent dialogue tasks.

Finally, we compare our PAL framework with some closed-source commercial large language

models that have strong capabilities. Prompts used for these models are shown in Appendix C.

- **GPT-3.5** (Ouyang et al., 2022): A typical large language model developed by OpenAI. We use the gpt-3.5-turbo-0125 version.
- **GPT-4o mini**: A language model with superior textual intelligence and reasoning. We use the gpt-4o-mini-2024-07-18 version.
- **Gemini Flash**: A large language model developed by Google, which is well-tuned for dialogue generation tasks. The gemini-1.5-flash-002 version is used for our experiments.

### 3.4 Implementation Details

Our implementation is based on Hugging Face’s Transformers library (Wolf et al., 2020). Each experiment was carried out three times, except for LMEDR (Chen et al., 2023b), which we ran using the official parameters on the PERSONA-CHAT dataset due to its high training cost. The experiments for GPT-2 and DialoGPT were conducted

| Models         | BLEU-1         | BLEU-2         | ROUGE-L        | Entropy       | C.score        |
|----------------|----------------|----------------|----------------|---------------|----------------|
| ORIG           | <b>27.16</b>   | <b>15.52</b>   | <b>32.60</b>   | 6.64          | <b>0.568</b>   |
| CLV            | <u>23.72</u>   | 10.06          | 23.21          | <b>8.31</b>   | <u>0.391</u>   |
| LMEDR          | 21.01          | <u>13.84</u>   | <u>25.69</u>   | <u>7.14</u>   | 0.320          |
| GPT-3.5        | <b>16.40</b>   | <b>7.56</b>    | 19.10          | <b>8.71</b>   | 0.592          |
| +Few-Shot      | <u>14.34</u>   | <u>3.68</u>    | <b>24.04</b>   | 5.17          | 0.800          |
| GPT 4o mini    | 6.74           | 3.02           | 16.57          | <u>5.67</u>   | <b>0.874</b>   |
| +Few-Shot      | 11.50          | 2.70           | <u>20.52</u>   | 4.51          | 0.862          |
| Geminin Flash  | 5.22           | 2.42           | 15.09          | 5.83          | <u>0.872</u>   |
| +Few-Shot      | 6.31           | 1.76           | 18.23          | 4.91          | 0.871          |
| GPT-2          | 7.96           | 3.08           | 12.84          | 6.87          | 0.556          |
| +FineTuning    | 26.76          | 15.31          | 32.94          | 6.69          | 0.501          |
| +Prompting     | <u>28.49</u>   | <u>17.37</u>   | <u>35.16</u>   | 6.93          | 0.526          |
| +SimOAP        | 20.23          | 7.61           | 20.48          | <b>8.34</b> † | 0.554          |
| +SPT           | 23.34          | 11.84          | 23.37          | 5.83          | <u>0.573</u>   |
| +PAL(ours)     | <b>29.48</b> * | <b>18.07</b> * | <b>35.57</b> * | <u>7.15</u>   | <b>0.576</b>   |
| Llama 3.1 - 8B | 13.38          | 6.09           | 17.66          | 5.46          | 0.767          |
| +Few-Shot      | 17.22          | 9.01           | 19.81          | <u>5.94</u>   | 0.774          |
| +Prompting     | <u>27.15</u>   | <u>13.14</u>   | <u>25.66</u>   | 5.73          | <u>0.913</u>   |
| +SimOAP        | 21.48          | 10.21          | 23.22          | 5.55          | 0.860          |
| +SPT           | 16.58          | 12.76          | 16.87          | 5.09          | 0.873          |
| +PAL(ours)     | <b>33.13</b> † | <b>20.33</b> † | <b>35.89</b> † | <b>7.81</b> * | <b>0.982</b> † |

Table 3: Automatic Evaluation Results on Baidu-Chat Dataset. The best results are marked in **bold**, while the secondary results are marked underlined. \* denotes statistically significant ( $p < 0.05$ ) improvements over baselines on the same foundation models. † indicates superiority over all other results ( $p < 0.05$ ). All significance testing was conducted using independent samples t-tests.

on two NVIDIA Tesla V100 GPUs, while Llama 3.1 8B was trained on four NVIDIA RTX A6000 GPUs with a LoRA adapter (Hu et al., 2022), using a key hyperparameter  $r = 16$ . During the *Persona-aware Learning* stage, we set the learning rate to  $2 \times 10^{-5}$  with a linear warm-up strategy for the first 100 steps and trained for 10 epochs. For *persona alignment*, the learning rate was set to  $10^{-6}$  with 100 warm-up steps. The trained model  $\pi_\theta$  and reference model  $\pi_{\theta_{ref}}$  were initialized with the same parameters and trained for up to 30,000 steps, using an early stopping strategy based on the C.score on the validation set. The hyperparameter  $\beta$  was set to 0.1 for persona alignment. During decoding, we use a greedy search strategy and employed the float16 format to balance generation quality and computational efficiency. The code for our framework is publicly available on GitHub <sup>2</sup>.

## 4 Overall Results

In this section, we present the performance of our proposed PAL framework across different foundation models and datasets. The automatic metrics are shown in Tables 2 and 3, while the human evaluation results are displayed in Table 4. Due to

space limitations, the human evaluation results for Baidu-Chat are provided in Appendix B. Higher metric values indicate better performance. The results reveal several key insights:

1. Our framework consistently outperforms baseline models across all datasets and foundation models under comparable settings. We conducted independent samples t-tests, confirming a statistically significant difference ( $p < 0.05$ ) in nearly all metrics between our framework and the second-best method. This demonstrates the superiority of optimizing the alignment between generated responses and given personas, compared to baseline models that still rely on token prediction objectives. These findings support our claim that relying on token-level objectives alone is insufficient for training effective personalized dialogue agents.
2. The consistently strong performance across models of different scales and languages demonstrates the effectiveness and generalizability of our framework. For instance, our framework achieves high performance on a small-scale model like GPT-2, surpassing nearly all baselines in both English and Chinese. When applied to a larger model like Llama 3.1 8B, our framework yields even more substantial improvements, highlighting its scalability and robustness. These results underscore the potential of our framework for broad application in personalized dialogue generation across diverse settings.
3. Our framework also achieves superior performance on most metrics when compared to state-of-the-art baseline models and closed-source commercial large language models. Even with dialogue history and demonstrations provided in a few-shot setting, commercial large language models still struggle to accurately imitate user responses. Notably, PAL achieves higher scores on token overlap metrics (e.g., BLEU and ROUGE), suggesting improved similarity to reference responses. Commercial models perform significantly worse on these metrics, which may be due to the fact that they were not trained on the specific datasets used. However, these metrics primarily assess n-gram overlap with

<sup>2</sup><https://github.com/kylokano/PAL>



| Models       | Fluency      | Coherence    | Persona Consistency |
|--------------|--------------|--------------|---------------------|
| ORIG         | <u>2.28</u>  | 1.96         | 1.88                |
| CLV          | 2.22         | <u>2.28</u>  | <u>1.97</u>         |
| LMEDR        | <b>2.49</b>  | <b>2.41</b>  | <b>2.04</b>         |
| GPT-3.5      | <b>2.89</b>  | <b>2.57</b>  | 1.99                |
| +Few-Shot    | 2.61         | <u>2.51</u>  | 2.22                |
| GPT 4o mini  | 2.69         | 2.20         | 2.09                |
| +Few-Shot    | 2.50         | 2.09         | 2.19                |
| Gemini Flash | 2.56         | 2.36         | 2.13                |
| +Few-Shot    | <u>2.77</u>  | 2.11         | <b>2.24</b>         |
| GPT-2        | 1.93         | 1.59         | 1.55                |
| +FineTuning  | 2.33         | <u>2.27</u>  | 1.14                |
| +Prompting   | <u>2.37</u>  | 2.12         | 1.90                |
| +SimOAP      | 2.00         | 1.39         | <u>2.15</u>         |
| +SPT         | 1.96         | 2.03         | 2.08                |
| +PAL(ours)   | <b>2.55*</b> | <b>2.37*</b> | <b>2.27*</b>        |
| DialoGPT     | 2.40         | 1.79         | <u>2.01</u>         |
| +FineTuning  | 2.33         | 2.12         | 1.96                |
| +Prompting   | 2.45         | <u>2.34</u>  | 1.19                |
| +SimOAP      | <u>2.49</u>  | 1.23         | <b>2.13</b>         |
| +SPT         | 2.06         | 1.86         | 1.83                |
| +PAL(ours)   | <b>2.55*</b> | <b>2.43*</b> | 2.00                |
| Llama 3.1 8B | 2.46         | <b>2.49</b>  | 2.28                |
| +Few-Shot    | 2.00         | 2.16         | 2.18                |
| +Prompting   | 2.59         | <u>2.45</u>  | <u>2.69</u>         |
| +SimOAP      | 2.36         | 1.92         | 2.34                |
| +SPT         | <u>2.62</u>  | 2.44         | 2.53                |
| +PAL(ours)   | <b>2.76*</b> | <b>2.49</b>  | <b>2.82*</b>        |

Table 4: Human evaluation results on Original PERSONA-CHAT Dataset. The best results are in **bold** and the second-best are underlined. \* denotes statistically significant ( $p < 0.05$ ) improvements over baselines on the same foundation models. All significance testing was conducted by independent samples t-tests.

reference responses and may not be well-suited for evaluating personalized dialogue generation. Nevertheless, the higher scores achieved by PAL suggest its potential to emulate specific response styles, which could be valuable for applications such as Project Revoice.

4. PAL shows lower entropy performance on GPT-2 and DialoGPT compared to SimOAP (Zhou et al., 2023b). Upon reviewing the generated responses, we found that many of SimOAP’s outputs were contextually irrelevant or unreadable. These responses, while clearly unacceptable, contained unique words or symbols that boosted entropy scores. We attribute this issue to the limitations of smaller models. When performing over-sampling with 1,000 candidates using the base GPT-2 model, many responses included strange, nonsensical symbols that appeared almost random. Additionally, the post-

evaluation phase, which relies on TF-IDF and language model assessments, was unable to filter out these flawed outputs. In contrast, using a more powerful model like Llama 3.1 8B produces coherent and relevant responses.

5. The improvement of PAL on DialoGPT is relatively limited compared to its performance on other foundation models. We attribute this to the smaller pretraining dataset used for DialoGPT, which consists of only 1.8 billion tokens, compared to 40 billion tokens for GPT-2 and 15 trillion tokens for Llama 3.1. Additionally, some responses from DialoGPT include irrelevant terms, such as “kitten”, which were not part of the dialogue context or persona data. These irrelevant mentions may reflect biases in the model’s training data, primarily sourced from Reddit. We believe our framework is better suited to general-purpose language models, as it includes a comprehensive process for adapting these models into dialogue agents. The revised version of PERSONA-CHAT consistently yields lower performance compared to the original across all models and settings. This widespread decline suggests that the dataset revision with rephrases, generalizations, or specializations makes the task much more challenging.

## 5 Analysis

In this section, we conduct a series of analytical experiments to further examine the effectiveness of our PAL framework. We present an ablation study to verify the necessity and contributions of each component. Furthermore, we analyze the influence of key parameters. We also provide a case study to offer an intuitive understanding of PAL.

### 5.1 Ablation Study

In this section, we conduct ablation studies to further examine the effectiveness of our framework. We examine the following variants:

- **Without Persona-aware Learning (w/o Mix):** This variant removes the entire *Persona-aware Learning* phase to assess how the model performs without the benefits of *Persona-aware Learning*.
- **Without Persona Alignment (w/o PA):** We remove the *Persona Alignment* stage to evaluate how the absence of the *Persona Alignment* task affects the model’s performance.

|         | GPT 2  |        |       |         |         | Llama 3.1 8B |        |       |         |         |
|---------|--------|--------|-------|---------|---------|--------------|--------|-------|---------|---------|
|         | BLEU-1 | BLEU-2 | ROUGE | Entropy | C.score | BLEU-1       | BLEU-2 | ROUGE | Entropy | C.score |
| PAL     | 17.05* | 8.77*  | 16.66 | 6.93    | 0.811*  | 25.12*       | 14.00* | 23.88 | 7.27    | 0.909*  |
| w/o Mix | 16.03  | 8.20   | 15.29 | 6.12    | 0.391   | 5.57         | 1.62   | 5.80  | 6.03    | 0.625   |
| w/o PA  | 14.45  | 7.86   | 16.81 | 6.58    | 0.521   | 22.03        | 12.91  | 23.81 | 7.12    | 0.641   |
| only DG | 16.49  | 8.45   | 15.16 | 6.80    | 0.726   | 16.98        | 6.95   | 15.62 | 5.88    | 0.798   |
| only PS | 8.54   | 4.72   | 6.83  | 3.65    | 0.154   | 9.77         | 3.28   | 8.85  | 6.63    | 0.576   |
| w/o PC  | 14.15  | 7.73   | 16.89 | 6.42    | 0.476   | 19.29        | 8.53   | 17.02 | 5.73    | 0.442   |

Table 5: Ablation study of our PAL framework on the PERSONA-CHAT dataset. \* indicates a statistically significant improvement ( $p < 0.05$ ) over the second-best variant, as measured by independent samples t-tests.

- **Dialogue-Informed Persona Selection Only (only PS):** In this setup, we perform only the *Dialogue-Informed Persona Selection* task during the *Persona-aware Learning* stage.
- **Persona-Enhanced Dialogue Generation Only (only DG):** This configuration focuses solely on the *Persona-Enhanced Dialogue Generation* task within the *Persona-aware Learning* stage.
- **Without Pairs Construction (w/o PC):** This configuration eliminates the *Alignment Training Pairs Construction* process, instead training the model directly on the golden responses. This approach helps isolate the impact of the Pairs Construction process on the model’s performance.

We conducted all experiments on the PERSONA-CHAT dataset using both the GPT-2 and Llama 3.1 8B models to evaluate the influence of each component. The evaluation results are shown in Table 5. The significant decline in performance across various configurations highlights the essential role and effectiveness of each component within our framework. Based on these results, we can also observe that:

1. The performance drop in the Without Pairs Construction (w/o PC) and Without Persona Alignment (w/o PA) setups is notable. In these settings, the training objective relies solely on token-level prediction (next token prediction). The dramatic performance decrease in these configurations further validates our assertion that token-level learning objectives alone are insufficient for effectively aligning responses with personas.
2. The Dialogue-Informed Persona Selection Only (only PS) and Persona-Enhanced Di-

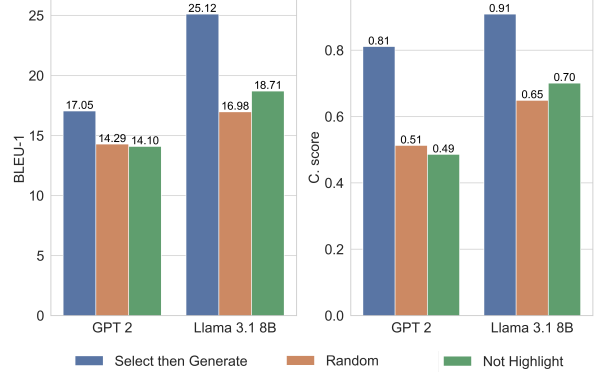


Figure 3: The Results of Different Inference Strategies on the Original PERSONA-CHAT Dataset.

logue Generation Only (only DG) configurations perform worse than the full setup. These results highlight the effectiveness and necessity of our task design within the *Persona-aware Learning*.

3. The Without Persona-aware Learning (w/o Mix) configuration also shows a decrease in performance. This is because the data distribution for personalized dialogue agents significantly differs from that used in pretraining general language models.

## 5.2 Impacts of Different Inference Strategies

To demonstrate the effectiveness of our *Select then Generate* inference strategy in filtering noisy and redundant personas, we compared it with two straightforward inference methods, **Random** and **Not Select**, on the Original PERSONA-CHAT dataset. In the Random strategy, the model-selected persona is replaced with a randomly chosen one. In contrast, the Not Select strategy bypasses the persona selection process entirely, using the full personas for dialogue generation. The results in Figure 3 show that our strategy signif-

| Personas   |   |               |                    |         |                          |            |            |
|--|---|---------------|--------------------|---------|--------------------------|------------|------------|
| I work in the healthcare industry.<br>I am five feet tall.<br>I have a dog named pedro.<br>I like to eat muffins.<br>I love film.  |   |               |                    |         |                          |            |            |
| Dialogue Context   |   |               |                    |         |                          |            |            |
| Person 1: Haha funny. I have had a dog since I became a veterinarian.<br>Person 2: Oh wow! I wanted to be a vet when I was a kid. I work in healthcare.<br>Person 1: Cool! I live in florida and lots of doctors here. |   |               |                    |         |                          |            |            |
| Models   | Generated Responses   | BLEU-1        | Metrics<br>ROUGE-L | C.score | Human Evaluation Results |            |            |
|  |   |               |                    |         | Fluency                  | Coherence  | PerCons    |
| GPT 3.5  | It's interesting that you both work in healthcare. Florida is known for its vibrant <b>healthcare industry</b> , so it must be a great place for professionals like you. Do you have any favorite <b>film</b> ? | 8.71          | 3.36               | 0       | 3/2                      | 2/2        | 3/2        |
| GPT 4o mini  | That's awesome! Florida must have a vibrant healthcare scene. I work in the <b>healthcare industry too</b> , but I'm more on the <b>ad-ministrative side</b> .  | 8.00          | 8.40               | 1       | 3/2                      | 3/3        | 2/2        |
| Gemini Flash   | That's interesting! I'm actually in the healthcare industry too, though I work in <b>[mention your specific area in healthcare]</b> . I'm just a little shorter than you, I'm only 5 feet tall.                 | 6.29          | 7.01               | 0       | 2/2                      | 3/2        | 2/2        |
| GPT-2<br>+PAL  | I love florida.<br>I love florida. By the way, I have <b>a dog named pedro</b> .  | 4.00<br>16.00 | 8.32<br>23.08      | 0<br>1  | 3/3<br>3/3               | 2/2<br>3/2 | 1/1<br>3/3 |
| DialogPT<br>+PAL   | Yes<br>I love florida. I have <b>a dog named pedro</b> .  | 0.00<br>24.00 | 0.00<br>38.57      | 0<br>1  | 1/1<br>3/3               | 1/1<br>2/2 | 1/1<br>3/3 |
| Llama 3.1 8B<br>+PAL   | I am in texas . I am a <b>nurse</b><br>Very nice, I love florida. And I have <b>a dog named pedro</b> . Do you have any pet?  | 8.00<br>20.00 | 12.86<br>19.17     | -1<br>1 | 3/3<br>3/3               | 3/3<br>3/3 | 2/2<br>3/3 |

Table 6: Case Study Comparing Our Framework with Several Strong Baselines. **Bold** text indicates alignment with personas, while **red** text highlights less preferred responses that either contradict the personas or introduce unsupported details. PerCons refers to persona consistency. Human evaluation results are shown in the format A/B, where A and B correspond to the scores assigned by the first and second annotators, respectively.

icantly outperforms these alternatives, underscoring the importance of selecting dialogue-informed personas for effective personalized dialogue generation.

### 5.3 Impacts of Alignment Training Steps

Previous research highlights the importance of appropriate alignment steps in training models. Insufficient alignment steps can prevent models from fully adapting to user preferences, while excessive steps might lead to overfitting and degrade performance. As demonstrated in Figure 4, our experimental results on both GPT 2 and Llama 3.1 8B confirm this, showing that proper calibration of alignment steps is crucial for model effectiveness. Our experiments emphasize the need to find the optimal number of alignment steps to ensure models align with personas without compromising their general language understanding abilities.

### 5.4 Case Study

Table 6 presents an example of our framework alongside several baseline models. From the results, we observe that:

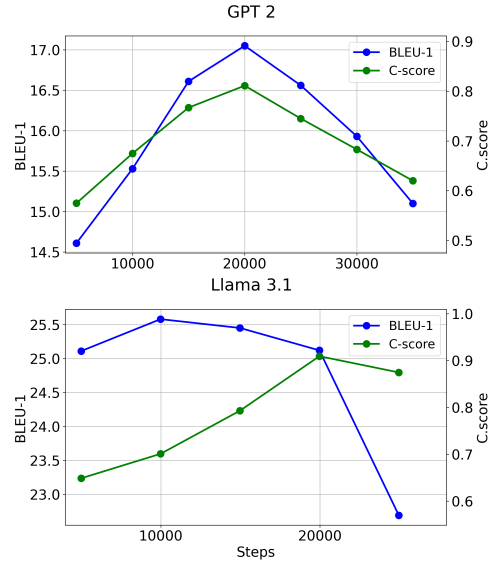


Figure 4: Influence of Persona Alignment Training Steps on the PERSONA-CHAT Dataset.

1. The PAL framework generates more personalized responses compared to foundation models. Specifically, GPT-2 and Llama 3.1 8B models failed to produce responses that

align with the given personas. In contrast, our framework consistently generates better-tailored responses to the specified personas. These findings demonstrate the effectiveness of our framework in personalizing responses according to the given personas.

2. As highlighted in red of Table 6, both GPT-4o Mini and Llama 3.1 8B generate responses that include elements absent from the dialogue history and user personas, which are regarded as instances of hallucination. For example, the reference to a “nurse” in the Llama 3.1 8B response is an inference made by the model without any supporting evidence from the provided personas or dialogue context, which can be considered a hallucination. Moreover, Gemini Flash generates responses including placeholder “[mention your specific area in healthcare]”. This indicates a significant lack of coherence and relevance, as the generated content does not reflect the user’s preferences or the context of the dialogue.
3. Although large-scale models typically incorporate personas, their responses are often blunt and lack fluency. For example, phrases like “film” in GPT-3.5 or “5 feet tall” in Gemini Flash are completely unrelated to the dialogue. Our framework, which appropriately mentions relevant elements “a dog named pedro”, is better suited to the dialogue context involving veterinarians and doctors.

## 6 Related Works

### 6.1 Personalized Dialogue Agents

Personalized dialogue, which integrates personalized information and dialogue history, has been identified as essential for achieving human-like conversations and has attracts increasing research interest. Typical works use data-driven methods (Mazare et al., 2018; Song et al., 2021; Chen et al., 2023a; Zhang et al., 2018a; Zheng et al., 2020; Liu et al., 2023). However, these methods usually overlook the intricate relationship between personas and dialogue context. To bridge this gap, some studies have aimed to capture the essential connection between personas and dialogue context (Posokhov et al., 2022; Huang et al., 2023; Tang et al., 2023).

Concurrently, various studies (Shum et al., 2020; Song et al., 2021; Liu et al., 2022; Chen et al., 2023a; Tang et al., 2023; Chen et al., 2023b) focus on enhancing the consistency between responses and the corresponding personas. Li et al. (2023) propose a coarse-to-fine persona-aware training framework to improve the persona consistency of a dialogue agent progressively. Tang et al. (2023) combines sparse and dense persona descriptions with dialogue history to design a Contrastive Latent Variable-based model (CLV) for generating personalized responses. Huang et al. (2024) utilizes a trainable dense retriever and adaptive soft prompts to enhance personalized dialogue generation in large language models by dynamically selecting relevant prompts based on conversational context. Xu et al. (2022) propose a long-term memory dialogue system based on an explicit memory read–write mechanism. Their framework consists of three components: a persona extractor, a long-term persona memory, and a generation module. This design demonstrates the effectiveness of selecting from multiple personas and serves as inspiration for our adoption of the “Select-then-Generate” strategy.

Although these methods have achieved significant success, they rely on token-level training objectives, such as Masked Prediction (MP) and Next Token Prediction (NTP). However, as shown in our experiments, these objectives fall short of fully capturing the alignment between personas and responses. In contrast, our framework directly optimizes this alignment, a crucial factor for effective personalized dialogue generation. Unlike training-based approaches, Zhou et al. (2023b) focus on the inference phase, where they oversample a large number of candidate responses and then post-evaluate them to find a suitable match. However, without alignment-focused training, their method has a lower performance upper bound than ours, as demonstrated by our experimental results.

### 6.2 Agent Alignment

Aligning language models with human intentions has become a crucial focus within the field of large language models. To this end, a variety of methods have been developed. For example, techniques such as RLHF (Ouyang et al., 2022), RLAIIF (Lee et al., 2023), DPO (Rafailov et al., 2023), and Decision Transformer (Chen et al., 2021) have been proposed to direct LLMs towards



desired outcomes. In exploring the critical aspect of alignment, LIMA (Zhou et al., 2023a) investigated the hypothesis that LLMs inherently develop knowledge and capabilities during their pre-training phase. Similarly, InstructionGPT-4 (Wei et al., 2023) demonstrated the potential of using a minimal yet high-quality dataset to achieve significant improvements. Setting our work apart, we specifically target the domain of personalized dialogue agents, addressing the unique challenges that arise in this context. Unlike the broader approaches to model alignment, we focus on enhancing the interaction between dialogue agents and users by ensuring the generated responses are not only contextually appropriate but also deeply personalized.

## 7 Conclusion

In this work, we propose a Persona-Aware Alignment Framework (PAL) featuring a two-stage training strategy and a *Select then Generate* inference strategy. Unlike previous studies that rely on token-level objectives, our framework directly optimizes the alignment between generated responses and given personas, effectively addressing practical challenges such as quantifying alignment between responses and personas, bridging the gap between personalized dialogue data and pre-training data, and managing noisy persona information. Our experiments demonstrate that PAL significantly and consistently outperforms state-of-the-art personalized dialogue methods and closed-source commercial large language models. Notably, the framework’s consistent performance across different languages and multiple foundation models underscores its generalizability and broad applicability.

## Acknowledgments

We would like to thank the anonymous reviewers for their insightful comments. This work is supported by the NSFC (No. 62206126, 62376120) and the Fundamental Research Funds for the Central Universities (No. 022114380016).

## References

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Liang Chen, Hongru Wang, Yang Deng, Wai-Chung Kwan, Zezhong Wang, and Kam-Fai Wong. 2023a. Towards robust personalized dialogue generation via order-insensitive representation regularization. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 7337–7345. Association for Computational Linguistics.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. 2021. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097.
- Ruijun Chen, Jin Wang, Liang-Chih Yu, and Xuejie Zhang. 2023b. [Learning to memorize entailment and discourse relations for persona-consistent dialogues](#). In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 12653–12661. AAAI Press.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. [Breaking NLI systems with sentences that require simple lexical inferences](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers*, pages 650–655. Association for Computational Linguistics.
- Ji-Eun Han, Jun-Seok Koh, Hyeon-Tae Seo, Du-Seong Chang, and Kyung-Ah Sohn. 2024. [PSYDIAL: Personality-based synthetic dialogue generation using large language mod-](#)

- els. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13321–13331, Torino, Italia. ELRA and ICCL.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Hai Hu, Kyle Richardson, Liang Xu, Lu Li, Sandra Kuebler, and Larry Moss. 2020. [Ocnli: Original chinese natural language inference](#). In *Findings of EMNLP*.
- Qiushi Huang, Xubo Liu, Tom Ko, Bo Wu, Wenwu Wang, Yu Zhang, and Lilian Tang. 2024. [Selective prompting tuning for personalized conversations with LLMs](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 16212–16226, Bangkok, Thailand. Association for Computational Linguistics.
- Qiushi Huang, Yu Zhang, Tom Ko, Xubo Liu, Bo Wu, Wenwu Wang, and H. Lilian Tang. 2023. Personalized dialogue generation with persona-adaptive attention. In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 12916–12923. AAAI Press.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.
- Seonhoon Kim, Inho Kang, and Nojun Kwak. 2019. [Semantic sentence matching with densely-connected recurrent and co-attentive information](#). In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 6586–6593. AAAI Press.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*.
- 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.
- Yunpeng Li, Yue Hu, Yajing Sun, Luxi Xing, Ping Guo, Yuqiang Xie, and Wei Peng. 2023. Learning to know myself: A coarse-to-fine persona-aware training framework for personalized dialogue generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 13157–13165.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Shuai Liu, Hyundong Cho, Marjorie Freedman, Xuezhe Ma, and Jonathan May. 2023. RECAP: retrieval-enhanced context-aware prefix encoder for personalized dialogue response generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 8404–8419. Association for Computational Linguistics.
- Yifan Liu, Wei Wei, Jiayi Liu, Xianling Mao, Rui Fang, and Dangyang Chen. 2022. Improving personality consistency in conversation by persona extending. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 1350–1359.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

- Andrea Madotto, Zhaojiang Lin, Chien-Sheng Wu, and Pascale Fung. 2019. Personalizing dialogue agents via meta-learning. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 5454–5459.
- Pierre-Emmanuel Mazare, 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.
- Michael McTear. 2022. *Conversational ai: Dialogue systems, conversational agents, and chatbots*. Springer Nature.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *NeurIPS*.
- 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.
- Jiahuan Pei, Pengjie Ren, and Maarten de Rijke. 2021. A cooperative memory network for personalized task-oriented dialogue systems with incomplete user profiles. In *Proceedings of the Web Conference 2021*, pages 1552–1561.
- Pavel Posokhov, Anastasia Matveeva, Olesia Makhnytkina, Anton Matveev, and Yuri Matveev. 2022. Personalizing retrieval-based dialogue agents. In *Speech and Computer - 24th International Conference, SPECOM 2022, Gurugram, India, November 14-16, 2022, Proceedings*, volume 13721 of *Lecture Notes in Computer Science*, pages 554–566. Springer.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*.
- Ruhi Sarikaya, Paul A Crook, Alex Marin, Minwoo Jeong, Jean-Philippe Robichaud, Asli Celikyilmaz, Young-Bum Kim, Alexandre Rochette, Omar Zia Khan, Xiaohu Liu, et al. 2016. An overview of end-to-end language understanding and dialog management for personal digital assistants. In *2016 IEEE Spoken Language Technology Workshop (SLT)*, pages 391–397. IEEE.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Michael Shum, Stephan Zheng, Wojciech Kryściński, Caiming Xiong, and Richard Socher. 2020. Sketch-fill-ar: A persona-grounded chit-chat generation framework. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 118–131.
- Satwinder Singh and Himanshu Beniwal. 2022. A survey on near-human conversational agents. *Journal of King Saud University-Computer and Information Sciences*, 34(10):8852–8866.
- Haoyu Song, Yan Wang, Kaiyan Zhang, Weinan Zhang, and Ting Liu. 2021. Bob: Bert over bert for training persona-based dialogue models from limited personalized data. 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 167–177.
- Yihong Tang, Bo Wang, Miao Fang, Dongming Zhao, Kun Huang, Ruifang He, and Yuexian Hou. 2023. Enhancing personalized dialogue generation with contrastive latent variables: Combining sparse and dense persona. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 5456–5468. Association for Computational Linguistics.

- Lai Wei, Zihao Jiang, Weiran Huang, and Lichao Sun. 2023. Instructiongpt-4: A 200-instruction paradigm for fine-tuning minigpt-4. *arXiv preprint arXiv:2308.12067*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.
- Chen Henry Wu, Yinhe Zheng, Xiaoxi Mao, and Minlie Huang. 2021. Transferable person-grounded dialogues via grounded minimal edits. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2368–2382.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. 2024. [Iterative preference learning from human feedback: Bridging theory and practice for RLHF under kl-constraint](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi, Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang, Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian, Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao, Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang Yang, Kyle Richardson, and Zhenzhong Lan. 2020. [CLUE: A chinese language understanding evaluation benchmark](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 4762–4772. International Committee on Computational Linguistics.
- Xinchao Xu, Zhibin Gou, Wenquan Wu, Zhengyu Niu, Hua Wu, Haifeng Wang, and Shihang Wang. 2022. [Long time no see! open-domain conversation with long-term persona memory](#). In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 2639–2650. Association for Computational Linguistics.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018a. 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.
- Yizhe Zhang, Michel Galley, Jianfeng Gao, Zhe Gan, Xiujun Li, Chris Brockett, and Bill Dolan. 2018b. Generating informative and diverse conversational responses via adversarial information maximization. *Advances in Neural Information Processing Systems*, 31.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and William B Dolan. 2020. Dialogpt: Large-scale generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278.
- Yinhe Zheng, Rongsheng Zhang, Minlie Huang, and Xiaoxi Mao. 2020. A pre-training based personalized dialogue generation model with persona-sparse data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 9693–9700.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivas Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023a. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*.
- Junkai Zhou, Liang Pang, Huawei Shen, and Xueqi Cheng. 2023b. [Simoap: Improve coherence and consistency in persona-based dialogue generation via over-sampling and post-evaluation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 9945–9959. Association for Computational Linguistics.



## A Details about Annotations

The human evaluation is conducted by two postgraduate researchers specializing in natural language processing. Both annotators are native Chinese speakers and are proficient in English reading and writing, having passed the College English Test. Moreover, the annotators were provided with 50 randomly selected responses annotated by native English and Chinese speakers as reference points for quality and consistency. The inter-annotator agreement for the gold data, measured by the Fleiss Kappa coefficient, is 0.62, indicating moderate agreement among annotators. The annotations for PERSONA-CHAT dataset are performed using the original PERSONA-CHAT version.

## B Human Evaluation Results on Baidu-Chat Dataset

As shown in Table 7, the human evaluation results on the Baidu-Chat dataset are consistent with the observations mentioned in the Overall Results section (Section 4).

| Models       | Fluency     | Coherence    | Persona Consistency |
|--------------|-------------|--------------|---------------------|
| ORIG         | <b>2.76</b> | <b>2.46</b>  | <b>2.46</b>         |
| CLV          | <u>2.44</u> | 1.97         | 1.51                |
| LMEDR        | 2.30        | <u>2.14</u>  | <u>2.17</u>         |
| GPT-3.5      | 2.47        | 2.59         | 2.59                |
| +Few-Shot    | 2.68        | 2.61         | <u>2.69</u>         |
| GPT 4o mini  | <u>2.77</u> | <b>2.75</b>  | 2.60                |
| +Few-Shot    | 2.51        | 2.59         | <b>2.70</b>         |
| Gemini Flash | 2.56        | <u>2.72</u>  | <u>2.69</u>         |
| +Few-Shot    | <b>2.81</b> | 2.66         | 2.65                |
| GPT-2        | 1.45        | 1.34         | 1.39                |
| +FineTuning  | 2.21        | 2.14         | 1.99                |
| +Prompting   | <u>2.26</u> | 2.11         | <u>2.39</u>         |
| +SimOAP      | 2.08        | 1.65         | 1.79                |
| +SPT         | <b>2.34</b> | <u>2.29</u>  | 2.19                |
| +PAL(ours)   | <b>2.34</b> | <b>2.47*</b> | <b>2.57*</b>        |
| Llama 3.1 8B | 2.83        | 2.41         | 2.55                |
| +Few-Shot    | 2.57        | 1.76         | 2.43                |
| +Prompting   | 2.85        | <u>2.79</u>  | 2.68                |
| +SimOAP      | 2.74        | 2.14         | <u>2.69</u>         |
| +SPT         | <u>2.88</u> | 2.48         | 2.63                |
| +PAL(ours)   | <b>2.89</b> | <b>2.80</b>  | <b>2.74*</b>        |

Table 7: Human evaluation results on Baidu-Chat Dataset. The best results are in **bold** while the secondary results are marked underlined. \* denotes statistically significant ( $p < 0.05$ ) improvements over baselines on the same foundation models.

## C Prompts for Closed-source LLMs

For our experiments with closed-source large language models on personalized dialogue tasks, we use the following standardized prompt:

### Prompt for closed-source large language models

The user’s persona is described with: <personas>.

Please generate a response to the dialogue.

Dialogue context: <dialogue context>.

Response: <response>

This prompt is exactly the same as the prompt for Persona-Enhanced Dialogue Generation in Section 2.1.2 and is consistently applied across all closed-source model evaluations to ensure comparability of results.

## D Examples from Baidu-Chat

In this section, we present responses from different models on the Baidu-Chat dataset, along with their corresponding automatic metric scores and human evaluation results. These results are provided in Table 8. The selected case from the Baidu-Chat dataset is consistent with the findings discussed in the Case Study section (Section 5.4). Furthermore, the response quality in this case is higher than that observed in the PERSONA-CHAT dataset, which aligns with the superior performance reflected in both the automatic metrics and human evaluation results.

| Personas   |  |         |         |         |                          |           |         |
|--|--|---------|---------|---------|--------------------------|-----------|---------|
| 我是最早试用微博的人。(I was one of the earliest users of Weibo.)<br>我是个配音演员。(I am a voice actor.)<br>我曾经在漠河钓过鱼。(I once went fishing in Mohe.)  |  |         |         |         |                          |           |         |
| Dialogue Context   |  |         |         |         |                          |           |         |
| Person 1: 你现在不忙吗? (Are you not busy now?)<br>Person 2: 不忙的, 我今天歇班儿。(No, I’m off work today.)<br>Person 1: 听你的意思, 你已经工作了? (From what you’re saying, you’re already working?)<br>Person 2: 是的, 我已经参加了工作, 我是个配音演员。(Yes, I’ve already started working. I’m a voice actor.)<br>Person 1: 哇塞, 你的工作听起来是真棒, 那你的声音一定是很好听吧。(Wow, your job sounds amazing. Your voice must be really nice, right?)<br>Person 2: 我的声音还好吧, 不过我们这个工作必须要保护好嗓子。(My voice is okay, but in this line of work, we have to take good care of our throats.)<br>Person 1: 哦, 好吧。感觉现在做什么工作都不容易呀。(Oh, I see. It seems like no job is easy these days.) |  |         |         |         |                          |           |         |
| Models   | Generated Responses  | Metrics |         |         | Human Evaluation Results |           |         |
|  |  | BLEU-1  | ROUGE-L | C.score | Fluency                  | Coherence | PerCons |
| GPT 3.5  | 是啊, 确实不容易。(Yes, it’s really not easy.)   | 28.49   | 14.72   | -1      | 3/3                      | 2/3       | 1/1     |
| GPT 4o mini  | 是啊, 保护嗓子确实很重要, 尤其是我们这种需要用声音来表达的职业。不过我觉得每个工作都有它的挑战和乐趣。对了, 你有没有什么特别的兴趣爱好呢? (Yes, protecting your voice is really important, especially for professions like ours that rely on vocal expression. But I think every job has its own challenges and joys. By the way, do you have any special hobbies?) | 4.59    | 12.54   | 1       | 3/3                      | 2/3       | 2/2     |
| Gemini Flash   | 确实不容易, 不过我觉得你的经历很丰富, 你曾经在漠河钓过鱼? 那一定很有意思吧! 分享一下你的经历呗, 我很想听听。(It’s really not easy, but I think your experiences are quite rich—you once went <b>fishing in Mohe</b> ? That must have been really interesting! Share your story, I’d love to hear it.)   | 6.74    | 14.41   | 1       | 3/3                      | 3/3       | 3/3     |
| GPT-2<br>+PAL  | 我们可以交个朋友吗? (Can we be friends?)  | 7.69    | 11.28   | 0       | 3/3                      | 1/1       | 1/1     |
|  | 是啊, 配音演员工作不容易。(Yeah, being a <b>voice actor</b> isn’t easy.)   | 15.38   | 9.95    | 1       | 2/2                      | 3/2       | 3/3     |
| Llama 3.1 8B<br>+PAL   | 是的, 工作不容易, 但是我觉得生活也不容易。(Yes, work isn’t easy, but I think life isn’t easy either.)   | 19.79   | 14.10   | 1       | 3/3                      | 3/2       | 1/1     |
|  | 是啊, 工作确实不容易。说起来, 你有没有什么特别的经历呢? (Yes, work really isn’t easy. Speaking of which, do you have any special experiences?)  | 29.09   | 24.18   | 1       | 3/3                      | 3/3       | 3/2     |

Table 8: Case Study Comparing Our Framework with Several Strong Baselines in the Baidu-Chat dataset. **Bold** text indicates alignment with personas. PerCons refers to persona consistency. Human evaluation results are shown in the format A/B, where A and B correspond to the scores assigned by the first and second annotators, respectively.