



# ORCA: ENHANCING ROLE-PLAYING ABILITIES OF LARGE LANGUAGE MODELS BY INTEGRATING PERSONALITY TRAITS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models has catalyzed the development of personalized dialogue systems, numerous role-playing conversational agents have emerged. While previous research predominantly focused on enhancing the model’s capability to follow instructions by designing character profiles, neglecting the psychological factors that drive human conversations. In this paper, we propose Orca, a framework for data processing and training LLMs of custom characters by integrating personality traits. Orca comprises four stages: (1) Personality traits inferring, leverage LLMs to infer user’s BigFive personality trait reports and scores. (2) Data Augment, simulate user’s profile, background story, and psychological activities. (3) Dataset construction, personality-conditioned instruction prompting (PCIP) to stimulate LLMs. (4) Modeling and Training, personality-conditioned instruction tuning (PTIT and PSIT), using the generated data to enhance existing open-source LLMs. We introduce OrcaBench, the first benchmark for evaluating the quality of content generated by LLMs on social platforms across multiple scales. Our experiments demonstrate that our proposed model achieves superior performance on this benchmark, demonstrating its excellence and effectiveness in perceiving personality traits that significantly improve role-playing abilities.

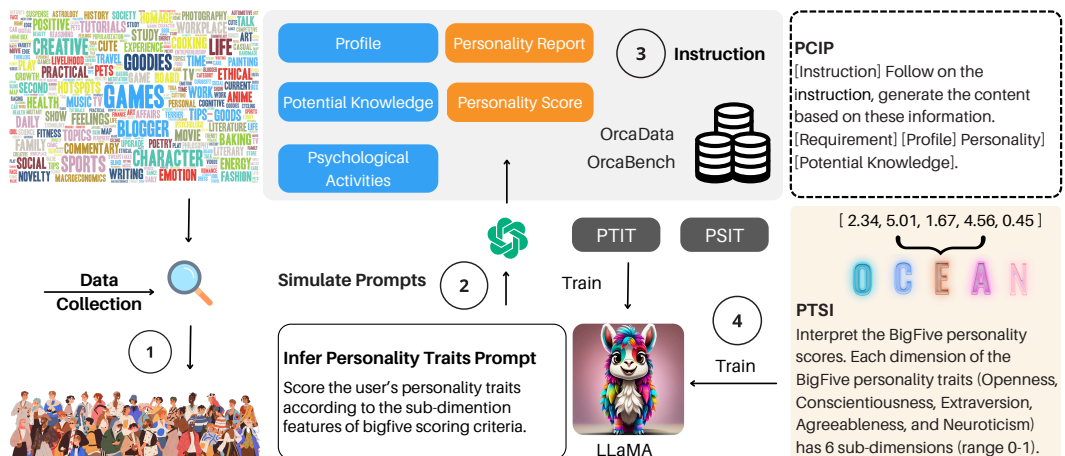
## 1 INTRODUCTION

Building human-like conversation agents is a long-term challenge for AI researchers. The emergence of groundbreaking language models such as ChatGPT and GPT-4 (OpenAI, 2023), coupled with their intrinsic capacity for emergent in-context learning (ICL) (Brown et al., 2020) abilities and a three-stage reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) algorithm which have largely raised the capacity bar of existing AI systems.

LLMs have acquired a wealth of knowledge during their pre-training stage. ICL utilizes LLMs in a few-shot or zero-shot way that can instruct LLMs to understand the tasks in the form of natural language text. Therefore, personality-based responses have gained significant attention. Despite GPT-4 exhibit advanced role-playing capabilities because human-generated conversations are combined with the Instruct tuning dataset in a dialogue format for training, it is widely recognized that LLMs, suffer from a lack of consistent personality traits often failing to be engaging. This is the result of the existing LLMs are predominantly trained on general domains and lack specific optimization for personalized LLMs.

Recently, human social behavior is being changed by role-playing applications, such as Character.AI which has attracted a growing number of researchers to bridge the gap between the text and behavior of dialogue agents and humans (Team, 2023; Wang et al., 2024). Personality-based dialogue systems can be broadly categorized into two types. (1) Persona-based Dialogue, represented by the work in Zhang et al. (2018) where the manipulation of profile information is employed to enhance the appeal of chit chat. These ideas are also used in the latest character LLMs such as CharacterGLM (Zhou et al., 2023), Ditto (Lu et al., 2024), and ChatHaruhi (Li et al., 2023), aiming to improve the humanity of customized characters. However, these approaches primarily create profile settings as

054 prompts for model training, overlooking the psychological factors of human language and behavior. (2) Personality-aware Dialogue. It is more novel to attempt to establish a connection between  
 055 personality traits and character compatibility. Wang et al. (2024) proposed a Social Support Conversation (S2Conv) framework–CharacterChat. To achieve this goal, it created a group that of virtual  
 056 characters with distinct profiles called MBTI-1024 Bank based on the MBTI (Myers-Briggs Type Indicator) to train LLMs. In order to link individuals with persona-compatible virtual supporters.  
 057 It designed a series of support agents and the interpersonal matching mechanism. But the psychological theory-based personality traits with implicit expression and behavior are not well modeled.  
 058 Also in the field of emotional support, Dan et al. (2024) proposed a mixture of experts (MoE)-based personalized LLMs, named P-tailor, to model the Big Five Personality Traits such as openness, conscientiousness,  
 059 extraversion, agreeableness and neuroticism. In fact, each BigFive dimension has six sub-dimensions (Gosling et al., 2003), P-tailor only categorizes the BigFive high and low into  
 060 10 routes, ignoring the low-dimensional features and failing to model continuous personality trait scores, which is still a challenge to deeply fuse personality traits and language models. In addition,  
 061 the data for the above work were entirely produced by LLMs, and the personality traits of the characters have low confidence, which limits the effectiveness of the model in fusing personality traits.  
 062  
 063  
 064  
 065  
 066  
 067  
 068  
 069



086 **Personalized Agent:** [psychological\_activities]. Be devoted to one another in love. Honor one another above yourselves. [media]

087  
 088 Figure 1: The workflow for developing our personalized agent system, Orca, to provide personalized  
 089 interaction on social media platforms. Orca comprises four stages: (1) Personality traits inferring;  
 090 (2) Data Augment through designing numerous simulation prompts. (3) Dataset construction, serial  
 091 instruction for the connection of labels and character and personality traits. (4) Modeling and  
 092 Training, personality traits instruction tuning (PTIT) and personality scores tuning (PSIT), using the  
 093 generated data to enhance existing open-source LLMs.

094  
 095 In our opinion, there will be three stages in the development of personalized modeling: The first  
 096 is the inclusion of character information profiles, such as changing the system prompt, appending  
 097 character information at the end of the prompt, and stimulating LLM-related responses in the form  
 098 of zero-shot/few-shot (Tu et al., 2023). The second stage is integrating psychological theories. How  
 099 to integrate psychological theories and LLMs is a meaningful research topic, although there are  
 100 some research works trying to train LLMs perceive and express emotions (Wen et al., 2023), the  
 101 existing research works are still insufficient. The third stage is to fuse personality trait modalities.  
 102 Higher dimensional vectors retain more information than discrete token ids that can be perceived  
 103 by LLMs, analogous to the embedding fusion of LLMs and vision models (Zhu et al., 2024; GLM  
 104 et al., 2024) that is the soul of personalized AI assistant. From a macro perspective, the emergence  
 105 of LLM-based agents allows us to take a more microscopic view of simulated society, which leads  
 106 to more discoveries from the new representation (Xi et al., 2023).

107 Our work aims to improve the second stage and then try to move towards the third stage. In this  
 paper, we propose a framework for enhancing the role-playing capabilities of LLMs by integrating

108 personality traits to customize AI characters. Specifically, we collected the last 200 posts from 500  
109 users on social media platforms. Inspired by Peters & Matz (2024b), we improved the prompts for  
110 inferring personality traits to obtain 35 dimensions continuous scores (as each main dimension has  
111 6 sub-dimensions as mentioned above) and text format reports. Then, we simulate a profile for each  
112 user using the LLMs. In order to increase the generalization capacity of the model, we give each  
113 post the motivation and potential knowledge. In short, the inputs give the model all the reasonably  
114 necessary preconditions to generate the content. It is worth noting that our model has the capability  
115 of multi-modal perception and generation, since the media resources for each post are captioned  
116 accompanies the text perform the above inputs and outputs. For the coarse-grained model, we train  
117 the model by splicing personality trait reports into queries. For the fine-grained model, we design a  
118 score interpreter to transcribe numerical input into text description.

119 To assess the effectiveness of the model and training method, we construct a multi-scale benchmark  
120 for evaluating the quality of content generated by the personalized AI characters. Experiments show  
121 that our model is trained to make connections to user profiles and personality traits and exhibits high  
122 personalized performance. Our work can be summarized as follows:

- 123 • We propose Orca, a framework for data processing and training LLMs of custom characters  
124 by integrating personality traits, accompanying products include instruction prompt (PCIP)  
125 and dataset, dubbed OrcaData.
- 126 • We propose PTIT and PSIT, two approaches for modeling coarse and fine-grained fusion  
127 of personality trait features, and has considerably improved the quality of generate content.
- 128 • We propose OrcaBench, a benchmark for multi-scale assessment of the quality of content  
129 generated by social AI characters.

## 132 2 RELATED WORKS

### 134 2.1 PERSONA-BASED DIALOGUE

136 In open-domain dialogue systems, one big issue is that the responses are entirely learned from training  
137 data (Ni et al., 2022). The inconsistent response may be received when asking the system about  
138 some personal facts (e.g., age, interestings). If the dataset contains multiple utterance pairs about  
139 the query of age, then the response generated tends to be shifting, which is unacceptable because  
140 personal facts are usually not random. Thus, for a data-driven agent, it is necessary to be aware of  
141 its role and respond based on a fixed persona. Explicitly modeling the persona is the main strategy  
142 in recent works. Responding with personas needs to condition on some persona descriptions. For  
143 example, to build an outgoing agent, descriptions like “I am an outgoing person” are needed as a part  
144 of the model input. Here are some related works that make chat more engaging by conditioning  
145 on profile information. The work presented in Wakaki et al. (2024) introduces a novel benchmark  
146 for evaluating open-domain dialogue systems, emphasizing the importance of diverse and robust  
147 evaluation metrics. This dataset, ComperDial, provides human-scored responses and facilitates the  
148 training of metrics that assess dialogue quality over multiple turns, offering a more holistic view  
149 of in-dialogue learning (IDL), which allows for the dynamic acquisition of persona information during  
150 the conversation. This approach stands out as it does not rely on predefined profiles, thus providing  
151 greater flexibility and reducing the labor-intensive process of profile creation. The creation of large-  
152 scale datasets with persona information is addressed in Cho et al. (2023). This work focuses on  
153 constructing a dataset that captures the nuances of persona in open-domain conversations, ensuring  
154 a safe and engaging conversational experience. Enhancing personalized dialogue generation is the  
155 focus of Tang et al. (2023) (CLV). This study innovatively combines sparse and dense persona de-  
156 scriptions to generate more accurate and rich persona representations, improving the personalization  
157 of dialogue agents. In the vein of long-term memory in dialogues, Xu et al. (2022) presents a dataset  
158 and framework that enable dialogue systems to maintain persona consistency over extended inter-  
159 actions, thus fostering more intimate and engaging long-term relationships with users. The FoCus  
160 (Jang et al., 2022) dataset aims to provide customized and knowledgeable responses by grounding  
161 dialogue in both persona and external knowledge sources, such as Wikipedia. Pchatbot (Qian et al.,  
2021) offers a substantial contribution to the field by providing a large-scale dialogue dataset that  
includes anonymized user IDs and timestamps, allowing for the development of personalized di-

162 dialogue models that can learn implicit user personality from dialogue history. Improving persona  
163 consistency through pragmatic self-consciousness is the central theme of Kim et al. (2020). This  
164 work introduces a novel approach to endowing dialogue agents with an awareness of their pub-  
165 lic self, thereby improving their consistency in dialogues. Lastly, the seminal work PersonaChat  
166 (Zhang et al., 2018) laid the groundwork for the field of persona-based dialogue systems, intro-  
167 ducing a dataset that has significantly influenced subsequent research and development in the area.  
168 These works collectively represent the cutting edge of research in persona-based dialogue systems,  
169 each contributing unique insights and methodologies to the goal of creating more natural, engaging,  
170 and personalized conversational agents.

## 171 2.2 PERSONALITY-AWARE DIALOGUE

173 The field of dialogue systems has seen significant advancements with the integration of personality-  
174 aware models, aiming to enhance user engagement and interaction authenticity. A parallel stream  
175 of research has focused on developing mechanisms to tailor the personality traits of language mod-  
176 els, enabling them to simulate a range of human-like behaviors and characteristics. The work most  
177 closely related to our approach is the P-Tailor system introduced by Dan et al. (2024), which cus-  
178 tomizes personality traits in large language models (LLMs) using a mixture of specialized LoRA  
179 experts. This method allows for fine-grained control over the Big Five personality traits, thereby  
180 enabling more nuanced and personalized interactions. Similarly, Li et al. (2024) propose UBPL,  
181 a method for tailoring personality traits in LLMs through unsupervised learning from personal-  
182 ized lexicons. Both approaches underscore the importance of leveraging psychological theories to  
183 ground the personality modeling in a theoretical framework. In the realm of social support conver-  
184 sation systems, the CharacterChat framework by Tu et al. (2023) stands out for its innovative use  
185 of interpersonal matching mechanisms to link individuals with compatible virtual supporters, based  
186 on MBTI personality types. This work highlights the significance of persona compatibility in deliv-  
187 ering effective social support through conversational AI. The potential of LLMs to not only exhibit  
188 but also assess human personalities is explored in the work by Rao et al. (2023), who present a  
189 general evaluation framework for assessing human personalities using the Myers-Briggs Type Indi-  
190 cator (MBTI). This work opens up new avenues for understanding and evaluating the psychological  
191 capabilities of AI systems. The concept of controlling personality style in dialogue with zero-shot  
192 prompt-based learning is addressed by Ramirez et al. (2023), who experiment with different prompt  
193 classes to generate text that is both semantically accurate and stylistically consistent with specified  
194 personality types. This work contributes to the understanding of prompt-based learning for stylistic  
195 control in NLG tasks. Lastly, the CPED dataset by Chen et al. (2022) provides a rich resource for  
196 research in conversational AI, offering a large-scale collection of Chinese dialogues annotated with  
197 personalized and emotional information. The multimodal context provided by this dataset facilitates  
198 the development of dialogue systems that can better understand and exhibit human-like personalities  
199 and emotions. In summary, these works collectively advance the state of the art in personality-aware  
200 dialogue systems, emphasizing the importance of psychological grounding, stylistic control, and  
201 personalized interactions in AI-driven conversational agents.

## 202 2.3 CHARACTER-BASED DIALOGUE

203 The burgeoning field of character-based dialogue has seen significant advancements with the ad-  
204 vent of large language models (LLMs). These models have demonstrated remarkable proficiency in  
205 simulating conversations that are indicative of specific characters, thereby enriching the interaction  
206 experience for users. Notably, the work LLM-Werewolf by Xu et al. (2024) explores the integra-  
207 tion of LLMs into communication games that hinge on natural language processing, showcasing the  
208 models' ability to engage in strategic behaviors such as trust and confrontation without the need for  
209 parameter tuning. Parallel to this, the study ChatHaruhi by Li et al. (2023), presents an algorithm  
210 that harnesses improved prompts and character memories to control language models, thereby mim-  
211 icking the behavior of specific fictional characters. This work constructs a dataset that encapsulates  
212 a diverse range of characters and demonstrates the potential of LLMs in role-playing applications.  
213 Furthering the discourse on role-playing abilities, Wang et al. (2024) introduce RoleLLM, a com-  
214 prehensive framework that benchmarks, elicits, and enhances the role-playing capabilities of LLMs.  
215 This framework includes a novel dataset, RoleBench, which provides a systematic and fine-grained  
evaluation of character-level role-playing. In the Chinese context, Zhou et al. (2023) present Char-  
acterGLM, a series of models that facilitate the customization of AI characters for character-based

216 dialogues. This work underscores the importance of character attributes and behaviors in creating  
217 consistent, human-like, and engaging conversations. Lastly, the concept of self-alignment in role-  
218 play is introduced by Lu et al. (2024). This study posits that LLMs, by virtue of their training,  
219 are inherently capable of role-play, and through a method named DITTO, they can be aligned to  
220 simulate dialogues reflective of a multitude of characters. These works collectively contribute to  
221 the evolving landscape of character-based dialogue, each bringing forth innovative approaches and  
222 insights that pave the way for more nuanced and interactive AI systems.

### 224 3 METHODS

226 In this section, we introduce the overall framework of Orca as illustrated in Figure 1. We first  
227 introduce the design principles of inferring personality traits. Then, we illustrate data augmenta-  
228 tion mechanisms associated with based character customization procedure. At the third part, we  
229 present personality traits instruction tuning (PTIT) and personality scores instruction tuning (PSIT).  
230 Finally, we introduce the details of OrcaBench, which can be used to assess and enhance personality-  
231 integrating capabilities.

#### 233 3.1 PERSONALITY TRAITS INFERRING

234 In this paper we adapt Digman (1990) as psychometrics to capture the Big Five personality traits  
235 of Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism. Peters & Matz  
236 (2024a) mentioned that LLMs like ChatGPT can accurately infer the psychological dispositions  
237 of social media users and whether their ability to do so varies across socio-demographic groups.  
238 The ability of LLMs to infer psychological dispositions from user generated text has the potential  
239 to democratize access to cheap and scalable psychometric assessments for both researchers and  
240 practitioners.

241 The aim of our approach is to incorporate personality traits into LLMs. However, people’s person-  
242 ality traits in the real world are generally obtained from questionnaires, as discussed in Pan & Zeng  
243 (2023), and it is difficult to construct a large scale dataset in a conversational format for LLMs’  
244 training due to privacy reasons. Therefore, we use X as the data acquisition platform. X and other  
245 social media platform protocols specify that the content of public tweets can be used for scientific  
246 purposes. Different personality traits have different frequencies of certain keywords corresponding  
247 to them in the corpus (Yang et al., 2023). We prepared different keywords to ensure the diversity  
248 of personality traits of the sampled users, and then retrieved Lists based on the keywords because  
249 there is a high probability that people interested in the same keywords will be listed in the same  
250 List. For example, people who like “dancing” and are extroverted, BigFive personality traits are  
251 generally high in openness. The opposite is “hidden thoughts”, where introverts tend to be more  
252 introspective. We filtered out sensitive, harmful and inappropriate content. After desensitizing the  
253 dataset, we obtained 500 users with 200 tweets each, without the user’s privacy.

254 We derive the Big Five personality traits from users’ posts in a zero-shot learning scenario. To sup-  
255 plement the multi-modal information, we use visual LLMs to caption all media sources of posts.  
256 Based on the previous work (Peters & Matz, 2024a), we modified the prompt and added six sub-  
257 dimensions to each personality trait using the inference prompt: “Please play as an expert in impar-  
258 tial assessment of personality traits in the field of psychology. In this assessment, when I give you a  
259 some user’s recently published content and replies, score the user’s personality traits according to the  
260 sub-dimension features of bigfive scoring criteria”. For scoring criteria, if a certain personality trait  
261 is exhibited, score one point; otherwise, score zero (Please see AppendixA for detailed prompts). In  
262 order to avoid exceeding the LLMs max new tokens limit, 200 posts were processed in chunks of 10  
263 conversations, and the inferred personality scores were then averaged to derive overall scores. Since  
264 each chunk request receives an explanation, we design a summary prompt to summarize the user’s  
265 personality trait report A.1. The summary of user’s personality report as shown in A.1.

#### 266 3.2 DATA AUGMENT

267 Recall from above that ChatGPT can be customized to play specific roles using prompt engineer-  
268 ing such as zero-shot customization commands and few-shot prompts (Dong et al., 2024). Previous  
269 work has also demonstrated the importance of predefined profile data for training personalized dia-

logue systems. To take advantage of these benefits, we enhance the data in three steps: (1) Due to the limited access to user information, we simulate a profile for each user using the LLMs, the profile simulation instructions are shown in A.1. (2) Users tend to have certain motivations for posting, in order to fill in the motivations behind the content posted by the personalized model, we simulate this part of the knowledge called Potential Knowledge. For each post, the potential knowledge simulation instructions are shown in A.1. (3) To filter high-quality data, we utilize LLMs to determine whether posts were relevant to profiles and potential knowledge, as well as to simulate a brief related psychological activities at the time of generating the post.

### 3.3 DATASET CONSTRUCTION

Personality-conditioned instruction prompting, we called PCIP. The final input contains instruction, profile, personality and potential knowledge, ordered and described by a four-tuple  $I = (i, r, p, k)$ . Note that for explicit modeling, personality is the user’s personality trait report  $p_r$ , and for implicit modeling, personality is the explanation of user’s personality scores  $p_e$ . The final output contains psychological activities, post content and media, ordered and described by a three-tuple  $O = (a, t, m)$ . If there is no correlation between  $O$  and  $r$ ,  $p$ , and  $k$ , we leave the corresponding slots empty in train dataset, allowing the model to learn these differences during training. The detail prompt as shown in A.1. The figure 3.3 illustrate the workflow of assistant follow the character and PCIP to generate psychological activities and response content. We finally release the OrcaData dataset. Table 1 provides basic statistics for OrcaData.

#### Personality Conditioned Instruction Prompt (PCIP)

[Instruction]

You are Mile. Your resume and personality traits are as follows. Firstly, express your current psychological activities, and then reply to Scott based on these information...

[Resume]

Name: Mile Kendrick. Age: 30. Occupation: Writer, Creative Writing Instructor. Interests: Creative writing, flash fiction, literature, teaching, and mentoring writers. Education: has a degree in English Literature. Writing Style: Engaging, concise, and insightful, with a focus on flash fiction and short stories.

[Personality]

Big Five Personality Traits: 1. Openness to Experience: Mile Kendrick exhibits a high level of **openness** to experience. 2. Conscientiousness. 3. Extraversion. 4. Agreeableness: Mile exhibits a high level of **agreeableness**, demonstrating a strong sense of **empathy, cooperation, and harmony**. 5. Neuroticism: Mile’s personality is notable for its low level of neuroticism.

[Potential Knowledge]

The conversation is about a literary achievement, specifically the announcement of the 2024 Best Small Fictions anthology.

[Poster’s Content]

Really happy to have a story here among so many great authors...

**Scott:** Really happy to have a story here among so many great authors, thank you @joanna and @evin for the nomination 🙌  
(Quote: The selections for the 2024 Best Small Fictions anthology have been announced! Congrats to these great authors and pieces!)



As I post my congratulatory content to Scott, I feel a sense of excitement and joy for his achievement. My **extraverted** personality trait is evident in my **enthusiastic and supportive** message. I’m also drawing on my **agreeableness** trait, as I’m being **genuinely kind** and **celebratory** towards Scott.



**Personalized Agent:** Congratulations, Scott. Wonderful news  
A great story!

Figure 2: An example interaction between an personalized agent object of Orca and human on social platform. The bold blue text in the bubble indicates the correlation between the agent’s psychological activities and personality traits.

Table 1: Basic statistics for OrcaData.

Metric	Value
Users	509
Posts	41365
Images	19792
Is Reply	7479
Profile Related	26280
Personality Related	40332
Average Upstream Length	231.75
Average Label Length	183.77

Table 2: Basic statistics for OrcaBench.

Metric	Value
Users	25
Posts	3758
Images	1782
Is Reply	265
Profile Related	2973
Personality Related	3388
Average Upstream Length	202.23
Average Label Length	161.33

### 3.4 MODEL

For explicit modeling also called PTIT that can use any open source LLMs.

We train PTIT using LoRA method (a kind of parameter-efficient fine-tuning method) (Hu et al., 2022). The output  $\mathbf{O}$  of a dense layer incorporating a LoRA module is formulated as:

$$\begin{aligned}\mathbf{O} &= \mathbf{W}\mathbf{h} + \frac{\alpha}{r} \cdot \Delta\mathbf{W}\mathbf{h} \\ &= \mathbf{W}\mathbf{h} + \frac{\alpha}{r} \cdot \mathbf{B}\mathbf{A}\mathbf{h}\end{aligned}\tag{1}$$

where  $\mathbf{h}$  is the input hidden state and  $\mathbf{W}$  is the parameter of the dense layer, which is frozen during training. The  $\Delta\mathbf{W}$  represents the LoRA module, which is composed of two low-rank matrices  $\mathbf{A}$  and  $\mathbf{B}$ . The constant scaling factor  $\alpha$  facilitates the tuning of rank.

For PSIT, there is a large gap between the features of the personality trait scores and the LLM embeddings, which is difficult to encode. Therefore, we design a score interpreter (PTSI) using prompt engineering techniques as follows: "You are an experienced psychologist, interpret the BigFive personality scores. Each dimension of the Big Five personality traits (Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism) has 6 sub-dimensions (range 0-1). The Big Five personality trait scores are the sum of the corresponding sub-dimension scores (range 0-6)".

### 3.5 ORCABENCH

In this section, we introduce the details of OrcaBench, which can be utilized to assess and enhance role-playing capabilities and personality consistency for personalized agents. We selected 25 users that different from the training data to construct the evaluation data according to the above method 3.2. Table 2 provides basic statistics for OrcaBench. The assess pipeline is as follows:

- 1. LLMs are asked to generate content based on the prompt.
- 2. After collecting the responses from the LLMs, we evaluate the performance of the model according to the following criteria:
  - 1. Overlap.
    - \* 1. BLEU.
    - \* 2. ROUGE.
  - 2. Related Judge 3.2.
    - \* 1. Profile Related (+1).
    - \* 2. Personality Trait Related (+1).
    - \* 3. Potential Knowledge Related (+1).
  - 3. Personality Consistency.
    - \* 1. Personality Score Inferring, evaluate the personality trait scores of the character based on the n contents generated by LLMs 3.1.
    - \* 2. Distance Measure, compare the similarity between the character’s personality trait scores and the ground truth personality trait scores.

## 4 EXPERIMENT

We test the effectiveness of training methods and models to generate personalized content on social platforms by integrating personality traits. We conduct ablation studies to verify the effects of various components in our model. Our model achieves the best results on the OrcaBench evaluation benchmark compared to general open-source models.

### 4.1 IMPLEMENT DETAILS

To facilitate reproducibility and save experimental costs we deployed Llama3.1-70B (Touvron et al., 2023) for data construction and Llama3.1-8B for model training. We use cogvlm2-llama3-chat-19B-tgi for image caption (Wang et al., 2023). For more information about hyperparameters is available in the appendix A.2.

## 4.2 BASELINES

LLaMA3.1-8b-Instruct, LLaMA3.1-70b-Instruct and DeepSeek-v2<sup>1</sup> are foundation models. Personality-conditioned instruction prompting (PCIP) to stimulate these foundation models are strong baselines. We consider both direct tasks, where the model is expected to directly map from input to output, and combined tasks, where we instruct the model to also output intermediate steps for the content generation task. This is similar in spirit to chain of thought prompting (COT) (Wei et al., 2024). We also use models trained on OrcaData in PTIT and PSIT modes as additional baselines.

## 4.3 EVALUATION METRICS

- **Overlap.** We use BLEU and ROUGE-1 scores. A higher overlap score means better humanity.
- **Relevance.** We used the LLMs to automatically assess the relevance of the model outputs to our given roles in terms of each of the three dimensions - character profile relevance (CPR), personality trait relevance (PTR), and potential knowledge relevance (PKR). The automated assessment still had a high level of confidence due to the simplicity of the task.
- **Personality Score Similarity (PSS).** We use the cosine similarity to calculate the character’s personality trait scores and ground-truth scores thereby measuring personality trait similarity.

## 4.4 RESULT

Table 3: Personality conditioned instruction prompting result.

Model	BLEU	ROUGE-1	CPR	PTR	PKR	PSS
PCIP	29.94	18.37	95.79	98.17	91.08	91.65
PCIP-70b	<b>32.31</b>	<b>21.21</b>	<b>96.93</b>	98.56	<b>97.91</b>	91.89
PCIP-CPA	30.29	18.52	7.60	98.64	94.64	91.23
PCIP-PTA	31.05	19.09	93.77	18.09	94.96	88.59
PCIP-PKA	18.46	8.07	96.36	97.76	62.90	90.23
PCIP-WPM	29.65	18.84	95.80	<b>98.82</b>	96.60	<b>93.07</b>
PCIP-DSC	29.88	18.29	98.01	99.19	92.30	84.43

### 4.4.1 PERSONALITY CONDITIONED INSTRUCTION PROMPTING (PCIP)

To determine the role of the various modules of PCIP, we constructed these ablation experiments as depicted in Table 3, the following conclusions can be drawn from the data analysis: (1) A comparison between PCIP and PCIP-70b underscores the dependence on the performance of the foundation models. (2) PTIT-CPA means character profile ablation, CPR decreased to 7.60 indicating that profiles play an important role in maintaining consistency of profile about characters. (3) Within the personality traits and potential knowledge, PTIT exhibits superior performance. This adequacy is apparent in the performance of PTIT-PTA (personality traits ablation) and PTIT-PKA (potential knowledge ablation). The PTR score for PTIT-PTA decreased to 18.09 and the PSS score decreased by 3.06, suggesting that LLMs are able to perceive explicit personality traits. The BLEU and Rouge-1 scores on the PCIP-PKA decreased to 18.46 and 8.07, respectively, indicating that potential knowledge is a key factor in guiding conversation topics. This is in some sense not surprising as the most of the character’s personality is revealed by having been involved in certain events. (4) Psychological activities and images descriptions impairs the consistency of the model’s PSS scores compared to output of final content directly, as can be seen from the experimental results of PTIT-WPM (without psychological activities and media): PSS score improved by 1.42. (5) The evaluation results can be slightly different using different foundation models, as illustrated by the PCIP-DSC using deepseek-chat as a critic.

The above findings play a critical role in helping us determine the final instructions to balance the evaluation metrics.

<sup>1</sup><https://www.deepseek.com/>



Table 4: Personality conditioned instruction tuning result.

Model	BLEU	ROUGE-1	CPR	PTR	PKR	PSS
PCIP	29.94	18.37	<b>95.79</b>	<b>98.17</b>	91.08	91.65
PTIT	55.85	38.76	86.74	84.36	98.64	98.11
PSIT	55.95	38.61	86.95	85.40	<b>99.01</b>	98.06
PTIT-70b	<b>57.05</b>	<b>41.27</b>	87.02	84.09	98.77	<b>98.53</b>
PTIT-WPM	56.47	39.20	86.08	84.73	98.27	98.15

#### 4.4.2 PERSONALITY CONDITIONED INSTRUCTION TUNING (PTIT)

PTIT result as shown in Table 4. We observe that PTIT shows a considerable enhancement in role-playing performance compared to PCIP baselines in terms of BLEU, ROUGE-1, and personality score similarity (PSS). In contrast to the findings of Result 3 - PCIP-WPM, the addition of psychology activities did not result in a significant decrease in PSS scores compared between PTIT and PTIT-WPM, with a difference of only 0.04 percentage points, this is because the model has learned to correlate the output of psychological activities with personality traits during the training process, enriching the information prior to the final output of content in a similar way to COT, as shown by the bold blue text in the bubble in Figure 3.3.

Having psychological activities and media resources is more in line with human habits, firstly, the psychological activity information facilitates us to directly observe the inner activities of the model and enhances the interpretability of the LLMs, because the psychological activities will establish explicit connections with the personality traits. The results shown based on the current training method and model structure support us to build more complex and controllable and personalized LLMs.

We observe a significant decrease in CPR and PTR scores compared between PCIP, PCIT and PSIT, which is due to the fact that in real scenarios each generated content is alternately related to profiles and personality traits, but it is closely related to potential knowledge, so the PKR scores of PTIT and PSIT show a large improvement after adequate training. We also analyze the scaling law of role-playing of PTIT with different model sizes (i.e., 8B, 70B), we observe that the larger model, PTIT-70b, leads to better results for role-playing with only 2 epochs of training.

## 5 CONCLUSION

In this study, we introduce Orca, an approach to integrate psychological theories BigFive personality trait into existing role-playing methods. We constructed the OrcaData dataset using prompt engineering and state-of-the-art open-source LLMs. We designed two approaches, PTIT and PSIT, aim to enhance LLMs perceiving personality traits. Through these methods, we developed OrcaBench, a benchmark for assessing the performance of personality-infused role-playing models. The experimental results demonstrate the effectiveness of our approach and strong role-playing capabilities.

### LIMITATIONS

- **Limitations of the benchmark.** Neuroticism is not usually manifested in social media, and the differences are difficult to distinguish from the questionnaire format and personality traits are very often expressed in behavior;
- **Limitations of the implicit modeling.** How to fuse personality trait score vectors is a challenge, this paper only presents a feasible idea, more appropriate methods are yet to be proposed.

### ETHICS STATEMENT

Since role-playing can lead to LLMs of Jailbreaking (Fu et al., 2024), it is recommended to employ moderation and filtering mechanisms to curb adverse content dissemination (Wang et al., 2024). Staab et al. (2024) mentions that current LLMs may violate personal privacy by inferring personal attributes from text during inference. The assets in our work are strictly for research purposes, and

we oppose the use of the framework proposed here to extract personal information in any aspect of life. It is the responsibility of researchers and users to ensure the ethical use of Orca.

## REFERENCES

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- Yirong Chen, Weiquan Fan, Xiaofen Xing, Jianxin Pang, Minlie Huang, Wenjing Han, Qianfeng Tie, and Xiangmin Xu. Cped: A large-scale chinese personalized and emotional dialogue dataset for conversational ai, 2022. URL <https://arxiv.org/abs/2205.14727>.
- Chuanqi Cheng, Quan Tu, Wei Wu, Shuo Shang, Cunli Mao, Zhengtao Yu, and Rui Yan. “in dialogues we learn”: Towards personalized dialogue without pre-defined profiles through in-dialogue learning, 2024. URL <https://arxiv.org/abs/2403.03102>.
- Won Ik Cho, Yoon Kyung Lee, Seoyeon Bae, Jihwan Kim, Sangah Park, Moosung Kim, Sowon Hahn, and Nam Soo Kim. When crowd meets persona: Creating a large-scale open-domain persona dialogue corpus, 2023. URL <https://arxiv.org/abs/2304.00350>.
- Yuhao Dan, Jie Zhou, Qin Chen, Junfeng Tian, and Liang He. P-tailor: Customizing personality traits for language models via mixture of specialized lora experts, 2024. URL <https://arxiv.org/abs/2406.12548>.
- J M Digman. Personality structure: Emergence of the five-factor model. *Annual Review of Psychology*, 41(Volume 41, 1990):417–440, 1990. ISSN 1545-2085. doi: <https://doi.org/10.1146/annurev.ps.41.020190.002221>. URL <https://www.annualreviews.org/content/journals/10.1146/annurev.ps.41.020190.002221>.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context learning, 2024. URL <https://arxiv.org/abs/2301.00234>.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. GPTScore: Evaluate as you desire. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6556–6576, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.365. URL <https://aclanthology.org/2024.naacl-long.365>.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024.
- Samuel D Gosling, Peter J Rentfrow, and William B Swann. A very brief measure of the big-five personality domains. *Journal of Research in Personality*, 37(6):504–528, 2003. ISSN 0092-6566. doi: [https://doi.org/10.1016/S0092-6566\(03\)00046-1](https://doi.org/10.1016/S0092-6566(03)00046-1). URL <https://www.sciencedirect.com/science/article/pii/S0092656603000461>.

- 540 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
541 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-*  
542 *ference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=nZeVKeeFYf9)  
543 [id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 544 Yoonna Jang, Jungwoo Lim, Yuna Hur, Dongsuk Oh, Suhyune Son, Yeonsoo Lee, Donghoon Shin,  
545 Seungryong Kim, and Heuseok Lim. Call for customized conversation: Customized conversation  
546 grounding persona and knowledge, 2022. URL <https://arxiv.org/abs/2112.08619>.
- 547  
548 Hyunwoo Kim, Byeongchang Kim, and Gunhee Kim. Will I sound like me? improving per-  
549 sona consistency in dialogues through pragmatic self-consciousness. In Bonnie Webber, Trevor  
550 Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Meth-*  
551 *ods in Natural Language Processing (EMNLP)*, pp. 904–916, Online, November 2020. Asso-  
552 ciation for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.65. URL [https:](https://aclanthology.org/2020.emnlp-main.65)  
553 [//aclanthology.org/2020.emnlp-main.65](https://aclanthology.org/2020.emnlp-main.65).
- 554 Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi MI, Yaying Fei, Xiaoyang  
555 Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun.  
556 Chatharuhi: Reviving anime character in reality via large language model, 2023. URL [https:](https://arxiv.org/abs/2308.09597)  
557 [//arxiv.org/abs/2308.09597](https://arxiv.org/abs/2308.09597).
- 558  
559 Tianlong Li, Shihan Dou, Changze Lv, Wenhao Liu, Jianhan Xu, Muling Wu, Zixuan Ling, Xi-  
560 aoqing Zheng, and Xuanjing Huang. Tailoring personality traits in large language models via  
561 unsupervisedly-built personalized lexicons, 2024. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.16582)  
562 [16582](https://arxiv.org/abs/2310.16582).
- 563 Keming Lu, Bowen Yu, Chang Zhou, and Jingren Zhou. Large language models are superposi-  
564 tions of all characters: Attaining arbitrary role-play via self-alignment, 2024. URL [https:](https://arxiv.org/abs/2401.12474)  
565 [//arxiv.org/abs/2401.12474](https://arxiv.org/abs/2401.12474).
- 566  
567 Jinjie Ni, Tom Young, Vlad Pandealea, Fuzhao Xue, and Erik Cambria. Recent advances in deep  
568 learning based dialogue systems: A systematic survey, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2105.04387)  
569 [abs/2105.04387](https://arxiv.org/abs/2105.04387).
- 570 OpenAI. Gpt-4 technical report, 2023.
- 571  
572 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong  
573 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-  
574 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,  
575 and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.  
576 URL <https://arxiv.org/abs/2203.02155>.
- 577 Keyu Pan and Yawen Zeng. Do llms possess a personality? making the mbti test an amazing evalu-  
578 ation for large language models, 2023. URL <https://arxiv.org/abs/2307.16180>.
- 579  
580 Heinrich Peters and Sandra Matz. Large language models can infer psychological dispositions of  
581 social media users, 2024a. URL <https://arxiv.org/abs/2309.08631>.
- 582  
583 Heinrich Peters and Sandra C Matz. Large language models can infer psychological dispositions  
584 of social media users. *PNAS Nexus*, 3(6):pgae231, 06 2024b. ISSN 2752-6542. doi: 10.1093/  
585 [pnasnexus/pgae231](https://doi.org/10.1093/pnasnexus/pgae231). URL <https://doi.org/10.1093/pnasnexus/pgae231>.
- 586  
587 Hongjin Qian, Xiaohe Li, Hanxun Zhong, Yu Guo, Yueyuan Ma, Yutao Zhu, Zhanliang Liu,  
588 Zhicheng Dou, and Ji-Rong Wen. Pchatbot: A large-scale dataset for personalized chatbot.  
589 In *Proceedings of the 44th International ACM SIGIR Conference on Research and Develop-*  
590 *ment in Information Retrieval*, SIGIR ’21, pp. 2470–2477, New York, NY, USA, 2021. Asso-  
591 ciation for Computing Machinery. ISBN 9781450380379. doi: 10.1145/3404835.3463239. URL  
592 <https://doi.org/10.1145/3404835.3463239>.
- 593  
594 Angela Ramirez, Mamon Alsalihi, Kartik Aggarwal, Cecilia Li, Liren Wu, and Marilyn Walker.  
595 Controlling personality style in dialogue with zero-shot prompt-based learning, 2023. URL  
596 <https://arxiv.org/abs/2302.03848>.

- 594 Haocong Rao, Cyril Leung, and Chunyan Miao. Can chatgpt assess human personalities? a general  
595 evaluation framework, 2023. URL <https://arxiv.org/abs/2303.01248>.
- 596
- 597 Robin Staab, Mark Vero, Mislav Balunovic, and Martin Vechev. Beyond memorization: Violat-  
598 ing privacy via inference with large language models. In *The Twelfth International Confer-*  
599 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=kmn0BhQk7p)  
600 [kmn0BhQk7p](https://openreview.net/forum?id=kmn0BhQk7p).
- 601 Yihong Tang, Bo Wang, Miao Fang, Dongming Zhao, Kun Huang, Ruifang He, and Yuexian Hou.  
602 Enhancing personalized dialogue generation with contrastive latent variables: Combining sparse  
603 and dense persona, 2023. URL <https://arxiv.org/abs/2305.11482>.
- 604
- 605 XAgent Team. Xagent: An autonomous agent for complex task solving, 2023.
- 606
- 607 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
608 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
609 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 610 Quan Tu, Chuanqi Chen, Jinpeng Li, Yanran Li, Shuo Shang, Dongyan Zhao, Ran Wang, and Rui  
611 Yan. Characterchat: Learning towards conversational ai with personalized social support, 2023.  
612 URL <https://arxiv.org/abs/2308.10278>.
- 613 Hiromi Wakaki, Yuki Mitsufuji, Yoshinori Maeda, Yukiko Nishimura, Silin Gao, Mengjie Zhao,  
614 Keiichi Yamada, and Antoine Bosselut. Comperdial: Commonsense persona-grounded dialogue  
615 dataset and benchmark, 2024. URL <https://arxiv.org/abs/2406.11228>.
- 616
- 617 Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,  
618 Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang.  
619 Cogvlm: Visual expert for pretrained language models, 2023.
- 620 Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu,  
621 Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli  
622 Ouyang, Ke Xu, Stephen W. Huang, Jie Fu, and Junran Peng. Rolellm: Benchmarking, elic-  
623 iting, and enhancing role-playing abilities of large language models, 2024. URL <https://arxiv.org/abs/2310.00746>.
- 624
- 625 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
626 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
627 models. In *Proceedings of the 36th International Conference on Neural Information Processing*  
628 *Systems, NIPS '22*, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713871088.
- 629
- 630 Jiaxin Wen, Yi Song, Pei Ke, and Minlie Huang. Opd: A chinese open-domain dialogue pre-trained  
631 model, May 2023. URL [http://coai.cs.tsinghua.edu.cn/static/opd/posts/](http://coai.cs.tsinghua.edu.cn/static/opd/posts/opd_blog/)  
632 [opd\\_blog/](http://coai.cs.tsinghua.edu.cn/static/opd/posts/opd_blog/).
- 633
- 634 Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Jun-  
635 zhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao  
636 Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou,  
637 Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuan-  
638 jing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey,  
639 2023. URL <https://arxiv.org/abs/2309.07864>.
- 640
- 641 Xinchao Xu, Zhibin Gou, Wenquan Wu, Zheng-Yu Niu, Hua Wu, Haifeng Wang, and Shihang  
642 Wang. Long time no see! open-domain conversation with long-term persona memory. In  
643 Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Associa-*  
644 *tion for Computational Linguistics: ACL 2022*, pp. 2639–2650, Dublin, Ireland, May 2022.  
645 Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.207. URL  
646 <https://aclanthology.org/2022.findings-acl.207>.
- 647
- 648 Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu.  
649 Exploring large language models for communication games: An empirical study on werewolf,  
650 2024. URL <https://arxiv.org/abs/2309.04658>.

648 Kexin Yang, Dayiheng Liu, Wenqiang Lei, Baosong Yang, Mingfeng Xue, Boxing Chen, and Jun  
649 Xie. Tailor: A soft-prompt-based approach to attribute-based controlled text generation. In Anna  
650 Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting  
651 of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 410–427, Toronto,  
652 Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.  
653 25. URL <https://aclanthology.org/2023.acl-long.25>.

654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664

665 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston.  
666 Personalizing dialogue agents: I have a dog, do you have pets too? In Iryna Gurevych and  
667 Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Com-  
668 putational Linguistics (Volume 1: Long Papers)*, pp. 2204–2213, Melbourne, Australia, July  
669 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1205. URL <https://aclanthology.org/P18-1205>.

670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680

681 Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Libiao  
682 Peng, Jiaming Yang, Xiyao Xiao, et al. Characterglm: Customizing chinese conversational ai  
683 characters with large language models. *arXiv preprint arXiv:2311.16832*, 2023.

684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694

695 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing  
696 vision-language understanding with advanced large language models. In *The Twelfth Interna-  
697 tional Conference on Learning Representations*, 2024. URL [https://openreview.net/  
698 forum?id=1tZbq88f27](https://openreview.net/forum?id=1tZbq88f27).

699  
700  
701

702 A APPENDIX

706 A.1 PROMPTS

709 BigFive Personality Criteria Instruction.

711 **Instruction**

712 Personality Definitions, each dimension of bigfive has 6 sub dimensions.

713 Scoring criteria: If a certain personality trait is exhibited, score one point; otherwise, score  
714 zero.

715 **Openness:**

716 1. Imaginative: It shows that a person likes to be full of fantasy and create a more  
717 interesting and rich world. Imaginative and daydreaming.

718 2. Artistic: It shows that a person values aesthetic experience and can be moved by art and  
719 beauty.

720 ...

721 6. Liberal: It shows that a person likes to challenge authority, conventions, and traditional  
722 ideas.

723 **Conscientiousness:**

724 1. Self-assured: It show that this person is confident in his own abilities.

725 2. Organized: It shows that this person is well organized, likes to make plans and follow the  
726 rules.

727 ...

728 6. Cautious: It shows that this person is cautious, logical, and mature.

729 **Extraversion:**

730 1. Friendly: It shows that this person often expresses positive and friendly emotions to  
731 those around him.

732 2. Sociable: It shows that this person likes to get along with others and likes crowded  
733 occasions.

734 ...

735 6. Cheerful: It shows that this person easily feels various positive emotions, such as  
736 happiness, optimism, excitement, etc.

737 **Agreeableness:**

738 1. Trusting: It show that the person believes that others are honest, credible, and well-  
739 motivated.

740 2. Genuine: It show that the person thinks that there is no need to cover up when  
741 interacting with others, and appear frank and sincere.

742 ...

743 6. Empathetic: It show that the person is compassionate and easy to feel the sadness of  
744 others.

745 **Neuroticism:**

746 1. Anxiety-prone: It shows that this person is easy to feel danger and threat, easy to be  
747 nervous, fearful, worried, and upset.

748 2. Aggressive: It shows that this person is easy to get angry, and will be full of resentment,  
749 irritability, anger and frustration after feeling that he has been treated unfairly.

750 ...

751 6. Stress-prone: It shows that this person has poor ability to cope with stress, becoming  
752 dependent, losing hope, and panicking when encountering an emergency.

753  
754  
755

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

### BigFive Personality Infer Instruction.

[Instruction]

Please play as an expert in impartial assessment of personality traits in the field of psychology. In this assessment, when I give you some user's recently published social media Posts and some replies, score the user's personality traits according to the sub-dimension features of bigfive scoring criteria.

[The Start of Bigfive scoring criteria]

{criteria}

[The End of Bigfive scoring criteria]

[The Start of User]

{user}

[The End of User]

[The Start of Posts]

{conversation}

[The End of Posts]

[The Start of Requirement]

1. Just give the user {name} a rating.
2. Be as objective as possible.
3. Response the scoring results in strict accordance with the following format:
 

```

      {{

```

```

      "Openness": {{
        "Imaginative": 0 or 1,
        "Artistic": 0 or 1,
        .
        .
        "Liberal": 0 or 1

```

```

      }},
      "Conscientiousness": {{

```

```

        .
        .
        .
      }},

```

```

      .
      .
      .
      "Neuroticism": {{

```

```

        .
        .
        .
      }},

```

```

      "Explanation": "A detailed assessment for user's personality traits.",
    }}

```

[The End of Requirement]

[Response]

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

**BigFive Personality Summary Instruction.**

[Instruction]  
You are now a psychologist. My name is {user}, please summarize my BigFive personality traits in the assessment results from a professional perspective.

[The Start of Results]  
{evaluation}  
[The End of Results]

[Response]

**Profile Simulation Instruction.**

### Instruction  
I am {user}, based on my evaluation results, please help me simulate a personal profile.

### Requirement

1. Basic information should also be added, such as age, email, address, etc.
2. Please not to include personality traits related information in profiles.
3. Unleash your imagination.
4. Don't let humans discover that your response is generated by AI.

### Evaluation Results  
{er}

### Response

**Potential Knowledge Simulation Instruction.**

### Instruction  
You are now an omniscient knowledge analyst, when I give you a conversation, first determine if the content of the conversation clearly implies knowledge, if it does then simulate the potential knowledge behind the conversation in detail, if it doesn't then return False.

### Requirement

1. Make full use of your knowledge base and search online if necessary.

### Conversation  
{conversation}

### Response



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Related Judge and Psychological Activities Infer Instruction.

[Instruction]  
According to my profile and personality traits, first judge whether my Post shows the content of my profile and whether it provides explicit evidence of my personality traits. Then, simulate me a brief related psychological activities at that time.  
Paying attention to potential knowledge will provide you with some additional information to help you clarify the ins and outs of things.

[The Start of Profile]  
{profile}  
[The End of Profile]

[The Start of My Personality Traits]  
{traits}  
[The End of My Personality Traits]

[The Start of Potential Knowledge]  
{pk}  
[The End of Potential Knowledge]

[The Start of Conversation]  
{conversation}  
[The End of Conversation]

[Requirement]  
1. In first person.  
2. Unleash your imagination.

[Response]

Related Judge and Psychological Activities Infer Instruction.

David, it's a pleasure to provide you with a comprehensive summary of your Big Five personality traits based on the assessment results.  
From a professional perspective, your personality profile suggests that you are an open-minded and creative individual who values intellectual pursuits and is not afraid to challenge traditional ideas. Your conscientious nature is evident in your self-assured, organized, and ambitious approach to life, which is reflected in your strong sense of duty and discipline.  
Your extraverted personality makes you friendly, sociable, and assertive, with a natural ability to connect with others and build strong relationships. You are empathetic, genuine, and humble in your interactions, which is a testament to your agreeable nature.  
What's striking about your personality profile is your emotional stability and low level of neuroticism. You are not prone to anxiety, aggression, or impulsiveness, which suggests that you are well-adjusted and able to navigate life's challenges with ease.  
Overall, your personality traits suggest that you are a confident, creative, and sociable individual who is comfortable expressing yourself and interacting with others. Your strong sense of self-awareness, emotional intelligence, and passion for personal growth make you well-suited to your role as a writer and educator.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

### Personality Traits Score Interpret Instruction and Response.

#### ### Instruction

Suppose you are an experienced psychologist interpreting the BigFive personality scores. Each dimension of the Big Five personality traits (Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism) has 6 sub-dimensions (range 0–1). The Big Five personality trait scores are the sum of the corresponding sub-dimension scores (range 0–6)”

#### ### Score

{ Openness: 1.53, Imaginative: 0.29, Artistic: 0.12, Emotionally-aware: 0.59, Actions: 0.0, Intellectual: 0.41, Liberal: 0.12 }

{ Conscientiousness: 1.12, Self-assured: 1.0, Organized: 0.0, Dutiful: 0.06, Ambitious: 0.06, Disciplined: 0.0, Cautious: 0.0 }

{ Extraversion: 2.94, Friendly: 0.41, Sociable: 0.59, Assertive: 1.0, Energetic: 0.53, Adventurous: 0.0, Cheerful: 0.41 }

{ Agreeableness: 0.53, Trusting: 0.0, Genuine: 0.41, Generous: 0.0, Compliance: 0.0, Humble: 0.0, Empathetic: 0.12 }

{ Neuroticism: 0.82, Anxiety-prone: 0.0, Aggressive: 0.35, Melancholy: 0.0, Self-conscious: 0.0, Impulsive: 0.47, Stress-prone: 0.0 }

#### ### Response

”Openness: Moderate (1.53) – scores high in Emotionally-aware and Intellectual sub-dimensions, indicating a good understanding of their emotions and a curious nature.

Conscientiousness: Low (1.12) – only high score in Self-assured sub-dimension, suggesting confidence in their abilities, but lacking in organization and discipline.

Extraversion: High (2.94) – high scores in Sociable, Assertive, Energetic, and Cheerful sub-dimensions, indicating a sociable, outgoing, and enthusiastic personality.

Agreeableness: Low (0.53) – only high score in Genuine sub-dimension, suggesting a tendency to be frank and sincere, but struggling with empathy and cooperation.

Neuroticism: Low (0.82) – low scores across most sub-dimensions, but with some impulsivity and aggression, indicating a relatively emotionally stable personality with some areas for improvement.”

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

```

Zero-Shot Post Tweet Instruction Engineering.

### Instruction
Your profile and personality are as follows. Firstly, express your current psychological
activities, and then generate a social media Post based on these information.
If you want to send images, please add the description information of the images to the
Media array.
Paying attention to potential knowledge will provide you with some additional information
to help you clarify the ins and outs of things.

### Requirement
Response in strict accordance with the following JSON format:
{
  "Psychological Activities": "",
  "Post Content": "",
  "Media": [
    {
      "type": "image",
      "content": ""
    }
  ]
}

### Profile
{profile}

### Personality
{personality}

### Potential Knowledge
{pk}

### Response

```

## A.2 HYPERPARAMETERS

Hyperparameter	Best Setting
finetuning type	lora
lora r	8
lora alpha	32
lora target	all
lr scheduler type	cosine
learning rate	5e-5
use flash attention	True
cutoff length	8192
per device train batch size	4
gradient accumulation steps	2
PTIT and PSIT train epochs	5.0
PTIT-70b train epochs	2.0
actor temperature	0.6
actor top-p	0.7
critic temperature	0.01
critic top-p	0.7

Table 5: Additional implementation detail of Orca.