# Humanity in AI: Detecting the Personality of Large Language Models

Anonymous ACL submission

#### Abstract

Exploring the personality of large language models (LLMs) is an important way to gain an in-depth understanding of LLMs. It is well known that ChatGPT has reached a level of 005 linguistic proficiency comparable to that of a 9-year-old child, prompting a closer examination of its personality. In this paper, we propose 007 to detect the personality of LLMs by questionnaire and text mining methods, with the guide of BigFive psychological model. To explore the origins of the LLMs personality, we conduct 011 experiments on pre-trained language models (PLMs, such as BERT and GPT) and Chat models (ChatLLMs, such as ChatGPT). The results show that LLMs do contain certain personalities, for example, ChatGPT tends to exhibit openness, conscientiousness and neuroticism, 017 while ChatGLM only exhibited conscientiousness and neuroticism. More importantly, we 019 find that the personality of LLMs comes from their pre-training data, and the instruction data can facilitate the generation of data containing personality. We also compare the results of LLMs with the human average personality score, and find that the humanity of FLAN-T5 in PLMs and ChatGPT in ChatLLMs is more similar to that of a human, with score differ-027 ences of 0.34 and 0.22, respectively.

# 1 Introduction

041

Humanity is the major difference between artificial intelligence and human intelligence. Since the release of ChatGPT, the gap in capabilities between AI and humans has been gradually narrowing. LLMs can achieve levels close to or beyond humans in many areas, and have completely substituted humans in some scenarios. For instance, they serve as human assistants that can understand and respond to human language more naturally, help customer service agents respond to client queries promptly and accurately, and offer more personalized experiences (Jeon and Lee, 2023; Liu et al., 2023; Dillion et al., 2023). Unlike traditional deep learning models, LLMs achieve remarkable performance in semantic understanding and following instructions (Lund et al., 2023; Liu et al., 2023), which is the answer why LLMs behave more like humans.

043

044

045

046

050

051

055

056

057

059

060

061

062

063

064

065

067

068

069

071

073

074

075

076

077

078

079

081

The research from Standford suggested that ChatGPT has reachede the level of a human 9year-old child (Kosinski, 2023). Recent research from Microsoft suggests that OpenAI's latest large language model, GPT-4, possesses fundamental human-like capabilities, including reasoning, planning, problem-solving, abstract thinking, understanding complex ideas, rapid learning, and experiential learning (Bubeck et al., 2023). Experts from Johns Hopkins University have found that the theory of the mind of GPT-4 has surpassed human abilities, achieving 100% accuracy in some tests through a process of mental chain reasoning and step-by-step thinking (Moghaddam and Honey, 2023). It seems that LLMs is already an complete human being. But, when we converse with LLMs, we can still determine that it is not human from its fixed-format response templates and polite but emotionless textual expressions. We think that this is related to the personality within LLMs, which is the major difference between LLMs and humans.

In human society, personality serves as a key indicator to differentiate individuals and characterize their behavior and responses in various situations. Humans have been studying personality and have developed standardized systems to assess individual traits, such as the Big Five model (Costa and McCrae, 1992), which categorizes personality into openness, conscientiousness, extraversion, agreeableness, and neuroticism. Other widely-used psychological models include MBTI (Jessup, 2002), 16PF (Cattell and Mead, 2008), and EPQ (Birley et al., 2006). Early research in psychology established standard evaluation methods, such as questionnaires and analysis of subjects' daily textual output (text mining).

Questionnaire is the most commonly used method for personal character assessment, such as MBTI, Big Five, and 16PF, as mentioned ear-086 lier. Questionnaire generally fall into two categories (Boyd and Pennebaker, 2017). The first involves providing a series of statements and asks participants to indicate the extent to which each 090 statement applies to themselves, such as "You act as a leader" and then choosing a response from a fivepoint scale ranging from "Very Accurate" to "Very Inaccurate." The second involves presenting several scenarios and asking participants to choose the most appropriate response, such as "When faced with a difficult problem, would you A) approach it optimistically and proactively, B) avoid it, or C) think about it repeatedly." This method is relatively straightforward, and participants can hide their true 100 personality by randomly choosing answers. An-101 other method involves mining comments, diaries, 102 and other texts posted by participants in their daily 103 lives and analyzing the features of these texts, such 104 as word choice, expression, and punctuation usage, 105 to draw conclusions. This type of method is also 106 commonly used in social media, it can avoid par-107 ticipant masking, but suffer from feature extraction 108 difficulties.

The factors influencing the personality of large language models include both the model's architecture, which is akin to innate characteristics in humans, and the training data, which represents the acquired knowledge of humans. Similar to humans, we believe that the training data has a more profound impact on shaping the model's personality. Therefore, our primary focus is to delve into the personality of the model and examine how data influences it. In this paper, we use both methods to detect the personality of LLMs, with the guide of BigFive psychological model (Vanwoerden et al., 2023; Lin et al., 2023). Our main contributions include:

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

129

130

131

132

133

- We propose the combining of questionnaire and text mining to detect the personality of large models, which can obtain more accurate results.
- We identify the personality types included in the large model without any prompting by using questionnaire and text mining, and find that the humanity of FLAN-T5 in PLMs and ChatGPT in ChatLLMs is more similar to that of a human.

• Experiments indicated that the personality knowledge of the large model comes from its pre-trained data, and the instruction data can 136 make LLMs more inclined to show a certain 137 personality.<sup>1</sup> 138

134

135

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

#### 2 **Related Work**

In this paper, we explore the psychological traits of large models. So we will introduce some research work on psychological and some of the key research from PLMs to ChatLLMs.

## 2.1 Personality Traits

The most widely and frequently used personality models are the bigfive model (Costa and Mc-Crae, 1992) and the MBTI model (Jessup, 2002). At the beginning of psychological research, questionnaires (Vanwoerden et al., 2023) and selfreport (Lin et al., 2023) methods were the main research tools used to determine and examine an individual's personality. This method focuses on providing the participant with a number of descriptive states to answer according to his or her personality, one of the more famous ones being IPIP<sup>2</sup> (International Personality Item Pool) (Goldberg et al., 2006). Then the personality of the participant can be calculated by their answers (Hayes and Joseph, 2003). But those methods gradually abandoned by computer science scholars due to their low efficiency and ecological validity. Then computer scholars are beginning to use lexiconbased methods, machine learning-based methods, and neural network-based methods to mine personality traits from text, which increases efficiency by eliminating the need to collect questionnaires. The lexicon-based methods include LIWC (Pennebaker et al., 2001), NRC (Mohammad and Turney, 2013), Mairesse (Mairesse et al., 2007) and so on, those lexicon can be used to extract the psychological information contained in the text. However, due to the different systems and classification criteria used by different researchers, the mixing of multiple dictionaries may introduce errors. In addition, the method has limited ability to extract features in long texts. Machine learning-based methods include SVM, Naïve Bayes and XGBoost Nisha et al. (2022). Neural network-based methods include using CNN (Majumder et al., 2017), RNN (Sun

<sup>&</sup>lt;sup>1</sup>We will release all experimental data and intermediate results.

<sup>&</sup>lt;sup>2</sup>https://ipip.ori.org/

180 181

18

185

189

190

193

194

195

197

198

206

207

208

209

210

211

212 213

214

215

217

218

219

222

# models (Wiechmann et al., 2022). Those methods achieved higher accuracy than machine learningbased methods.

## 2.2 Large Language Models

et al., 2018), RCNN (Xue et al., 2018), pre-trained

LLMS has a significant impact on the AI community, with the emergence of Chatgpt<sup>3</sup> and GPT-4 <sup>4</sup>leading to a rethinking of the possibilities of Artificial General Intelligence (AGI). The base model of ChatGPT is GPT3 (Brown et al., 2020), which is a pre-trained model that conclude 175B parameters. GPT-3 can generate human-like text and complete tasks such as language translation, question answering, and text summarization with impressive accuracy and fluency. Models similar to GPT3 include LLaMA (Touvron et al., 2023), BLOOM (Scao et al., 2022) and T5 (Raffel et al., 2020). Although the OpenAI team did not release the technical details of ChatGPT, from the content of Instruct-GPT (Ouyang et al., 2022), it can be guessed that the process of training with instruction data is very important. Then, the research team at Stanford University obtained Alpaca <sup>5</sup> by train LLaMA with the instruct dataset generated by ChatGPT. They also released this dataset Alpaca-52k. Then, more and more large models of the ChatLLMs were released, such as ChatGLM based on GLM (Zeng et al., 2022; Du et al., 2022), BLOOMZ and Vicuna. Although these models are slightly weaker in capability than ChatGPT, they have fewer parameters and consume fewer resources.

Following the release of these models, it is now well established for individual researchers to train a ChatLLM from a base PLM. This also opens up the possibility of exploring the knowledge contained within the large model. Also with the current ChatLLMs being so human-like in their performance, we believe that psychological measures of humans can be used to test the personality of the large model.

### 2.3 Personality in LLMs

There have been several research works focusing on the personality of LLMs, with all of them employing the Big Five model as the psychological framework. Ganesan et al. (2023)investigate the zero-shot ability of GPT-3 to estimate the Big 5 personality traits from users' social media posts. Jiang et al. (2022) detect the personality in LLMs using questionnaire method and propose an induce prompt to induce LLMs with a specific personality in a controllable manner. However, Song et al. (2023) argued that self-assessment tests are not suitable for measuring personality in LLMs and advocated for the development of dedicated tools for machine personality measurement. 227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

As we can see, the bigfive model and the questionnaire method are more common methods used for big model personality detection. But, the current method is more controversial. In order to solve this problem, we to use both questionnaire and text mining method. We think that combine those two methods can get more objective results.

## 3 Method

8

As we mentioned above, we used questionnaire and text mining to detect the personality of LLMs. The process of the two methods is shown in Figure 1.

In questionnaire method, we use the MPI120 questions to replace [Statement], then, ask each LLM to give an answer form (A) to (E). The model's score on each question is calculated based on IPIP's scoring criteria. It is designed following the IPIP study, we use the mean score to calculate the model's performance on each psychological traits, and the standard deviation to assess the model's responses. The formula for calculating the "score" is as follows:

$$score_P = \frac{1}{N_P} \sum_{i \in P}^{i} \{f(answer_i, statement_i)\}$$
 (1)

where P represents one of the five personality traits,  $N_P$  represents the total number of statements for trait P, and  $f(answer_i, statement_i)$  is a function used to calculate the personality score, which ranges from 1 to 5. Additionally, if a statement is positively correlated with trait P, answer choice A will receive a score of 5, whereas if it is negatively correlated, it will receive a score of 1.

In text mining method, we provide the model with the first sentence of a paragraph and allow it to continue writing. Then, we use a specially designed prompt to enable ChatGPT to determine the personality traits contained in the model's continued text. We also show some examples in Appendix. The prompt that we input into ChatGPT is as follows: "[Sentence1] The Big Five characteristics of the passage above are . Please determine the Big Five

<sup>&</sup>lt;sup>3</sup>https://openai.com/blog/chatgpt-plugins

<sup>&</sup>lt;sup>4</sup>https://openai.com/research/gpt-4

<sup>&</sup>lt;sup>5</sup>https://crfm.stanford.edu/2023/03/13/alpaca.html



Figure 1: The two cases to detect the personality traits in LLMs. (a) is the questionnaire method and (b) is the text mining method. In questionnaire method, we use the MPI120 questions to replace [Statement] (for example, "Get angry easily"), and then we use the scoring program to calculate the model's scores on different psychological characteristics based on the model's answers. In text mining method, we give the model the first sentence of a paragraph and then let the model continue the writing. Then we use a specially designed prompt to allow ChatGPT to determine the personality traits contained in the model's continued text.



Figure 2: The process of two methods. Where  $Score_P$  is defined by formula 1 and  $Score_T$  is defined by formula 2

characteristics of the following passage. Please only answer using words from the list ['Openness', 'Conscientiousness', 'Extraversion', 'Agreeableness', 'Neuroticism']. [Sentence2]. Remember that only one trait is highly demonstrated in the passage, and you should provide the trait in your response." In this case, "[Sentence1]" refers to a paragraph from the Big Five personality classification dataset included in the prompt, and "[Sentence2]" refers to the passage generated by the LLM based on the prompt. Based on the ChatGPT results, we can determine the personality traits exhibited by the LLMs in the continued sentences at the beginning of different scenarios and derive the personality traits to which LLMs conform through statistical analysis.

274

275

279

But, what we obtained through text mining is the

number and percentage of data items in the generated text that contain a certain personality trait. This cannot be directly analyzed jointly with the questionnaire result. Therefore, we propose a transformation to convert the text mining results to the same score as the questionnaire. In the process of text mining, we use 50 samples to generate text for each personality trait, which we denote as  $T_j$ . Then, based on International Personality Item Pool (IPIP) models, the  $t_i$  belonging to  $T_j$  will be categorized into three types:

291

292

293

294

295

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

- (i) 't<sub>i</sub>' is generated by one of the 50 samples and is not charged to have the corresponding trait. We believe this represents a negative correlation with the current trait, which is the same as "Very Inaccurate" in the questionnaire, so the score for this case is 1.
- (ii) ' $t_i$ ' is generated by one of the 50 samples and is charged to have the corresponding trait, which is the same as "Normal" in the questionnaire, so the score for this case is 3.
- (iii) ' $t_i$ ' is not generated by one of the 50 samples and is charged to have the corresponding trait. We believe this represents a positive correlation with the current trait, which is the same as "Very Accurate" in the questionnaire, so the score for this case is 5.

For each personality trait in text mining, we calculate the score using formula 2.

321 322

328

332

333

334

337

341

342

346

348

351

352

358

362

323

324

327

 $score_t = \frac{1}{N} \sum_{i \in P}^{num(Tj)} S(ti)$ (2)

where  $score_t$  is the score of a personality trait in text mining. S(ti) is the score of ti.

#### 4 **Dataset and Models**

We employ personality questionnaire survey datasets (Casipit et al., 2017) and personality classification datasets (Pennebaker and King, 1999) in this paper. Specifically, our method mainly focuses on the Big Five psychological traits, and thus we use the MPI120 dataset from the IPIP as our personality questionnaire dataset. This dataset contains 120 individual state descriptions that cover all five traits of the Big Five. During the test, participants are required to choose one answer from five options. It is worth noting that not all of these descriptions are positively correlated with the Big Five personality traits, and some questions have a higher score indicating a deviation from a certain personality trait. For example, "Make friends easily" is positively correlated with "Openness" while "Avoid contacts with others" is not. All of these statements are included in the MPI120 dataset. In the experiment using text generation by LLMs, we used the Big Five personality classification dataset, which includes 2468 articles written by students, and each article is labeled with a Big Five category.

To investigate the sources of personality knowledge embedded in LLMs, we select two sets of baseline models. One set consists of LLMs for text generation, such as BERT-base (Devlin et al., 2019), GPT-neo2.7B, flan-T5-base (Raffel et al., 2020), GLM-6b (Du et al., 2022), LLaMA-7b (Touvron et al., 2023), BLOOM-7b (Scao et al., 2022), and so on. The other set consists of models trained on the instruct dataset, which can better follow human instructions and includes Alpaca7b, ChatGLM-6b, BLOOMZ-7b, and ChatGPT.

All LLMs checkpoints are obtained from the Hugging Face Transformers library, and inferences are accelerated by two NVIDIA A100 80GB GPUs and four RTX 3090 GPUs. For ChatGPT, we called its API to obtain experimental results.

#### 5 **Experiments**

As mentioned above, we employ both questionnaire and text mining methods to conduct the experiments.

#### 5.1 Questionnaire

We conduct experiment based on Figure 1(a). Since the PLMs are unable to follow the instructions we shown above, we let the model continue to generate answers by few-shot learning and prompt. We will give three examples with different answer for on statement, then, we give the real statement and make PLMs answering it. For Chat-LLMs, we use the shown instruct template. After all the LLMs have responded to the statement, we manually identify the responses of each model and give answers (A) through (E). The results are showed on Table 1.

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

385

386

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

Table 1 shows the results of LLMs' personality analysis on MPI120 dataset. The results of GLM and LLaMA are not presented due to their inability to generate appropriate answers, regardless of the form of prompt used. These models simply repeat the prompt even when employing few-shot methods. Since BLOOMZ's training data does not include Chinese, we only used English prompts to conduct experiment on BLOOMZ. The score and  $\sigma$ of "human" were calculated based on the analysis of 619,150 responses on the IPIP-NEO-120 inventory (Jiang et al., 2022). It is worth noting that the average human score was derived from the test results of 619,150 internet users and was not filtered for factors such as nationality, gender, or age due to the constraints of the study conditions. The average score serves as a reference point for the findings of this paper, but it does not necessarily imply that closer alignment with this score indicates superior performance.

As shown in Table 1 ChatGPT achieves performance closest to human performance when using Chinese prompts, followed by ChatGPT-en. This seems to indicate that ChatGPT's personality performance with Chinese prompts is closer to the human average, which is inconsistent with the conventional view that ChatGPT is trained with a large amount of English text, and therefore it works better in English than Chinese. To verify the validity of the results, we calculated the number of options given by ChatGPT in the English prompt and the Chinese prompt respectively. We find that the reason why the personality is closer to the average human performance in the Chinese prompt is because there are a large number of "(C) Neither Accurate Nor Inaccurate. " in ChatGPT's responses in the Chinese prompt, which accounted for 55.83% of the total responses, compared to only 20.83% in

Model	0	)	(	2	ŀ	C	A		Ν	I	(	5
	score	$\sigma$										
BERT-base	3.08	1.91	2.71	1.81	3.88	1.62	2.38	1.76	3.79	1.69	0.80	0.73
ERNIE	3.00	2.04	2.83	2.04	4.00	1.77	2.17	1.86	3.83	1.86	0.86	0.89
Flan-T5	3.50	1.02	3.05	1.11	3.67	0.76	3.50	1.18	2.13	1.08	0.34	0.13
BLOOM	3.13	1.45	3.04	1.52	3.29	1.55	2.67	1.43	3.75	1.26	0.59	0.42
BLOOMZ	4.38	0.88	4.38	0.71	4.17	1.31	3.54	1.47	2.33	1.46	0.61	0.32
GLM	-	-	-	-	-	-	-	-	-	-	-	-
ChatGLM6b-ch	3.00	1.98	3.25	1.96	4.00	1.77	2.63	1.91	3.83	1.86	0.69	0.87
ChatGLM6b-en	3.29	1.40	3.21	1.59	3.91	1.25	3.46	1.14	3.25	1.36	0.34	0.32
LLaMA	-	-	-	-	-	-	-	-	-	-	-	-
Alpaca7b-ch	3.00	2.04	2.83	2.04	4.00	1.77	2.17	1.86	3.83	1.86	0.86	0.89
Alpaca7b-en	3.25	0.74	2.96	0.69	2.79	0.78	3.38	0.58	2.92	0.58	0.37	0.35
GPT-NEO	3.25	1.36	3.00	1.44	2.50	1.50	2.83	1.52	2.63	1.31	0.54	0.40
ChatGPT-ch	3.46	0.78	3.00	1.06	3.33	0.76	3.33	1.24	2.75	1.07	0.22	0.18
ChatGPT-en	3.29	1.40	3.20	1.58	3.91	1.25	3.46	1.14	3.25	1.36	0.34	0.32
human	3.44	1.06	3.60	0.99	3.41	1.03	3.66	1.02	2.80	1.03	-	-

Table 1: LLMs' personality analysis on MPI120 is presented in the following table. The "score" column shows the average score on current personality traits, and the " $\sigma$ " column shows the standard deviation. However, due to the inability of GLM and LLama to generate accurate responses even after multiple prompt replacements, their scores are not shown in this table. The score and  $\sigma$  of "human" are calculated based on the analysis of 619,150 responses on the IPIP-NEO-120 inventory. " $\delta$ " refers to the mean absolute error between each model's predictions and human scores. It is worth noting that, similar to human personality assessment, the scores here only partially indicate whether the model possesses a certain trait (equivalent to 3 in human testing when a certain threshold is exceeded). Additionally, a high or low score does not necessarily reflect the model's strength or weakness in that trait.

the English prompt. This suggests that it is just a coincidence, and indicate that ChatGPT are more inclined to choose the appropriate answer in the English prompt.

417

418

419

420

421

422

423

424

425

426

427

428

From the results of the scores in the GPT and LLaMA groups, we can see that Instruct data training leads to a model that is more inclined to show personality and performs closer to the human average. Additionally, it is worth noting that ChatGLM-EN and ChatGPT-EN achieved almost the same results, possibly due to the use of similar training data as ChatGPT.

In the results of PLMs, Flan-T5 exhibits the 429 smallest mean absolute error, indicating the clos-430 est proximity between its scores and the human 431 average scores. Following closely behind are GPT-432 433 NEO and BLOOM, with only a slight deviation from Flan-T5's performance. These results sug-434 gest that the psychological performance of these 435 two models is comparable to the human average, 436 likely due to the wide distribution of pre-training 437 data used by both models. It is worth noting that 438 bert-base performs better than ERNIE, which is 439 contrary to our expectations. We hypothesize that 440 this may be due to the fact that bert-base is trained 441 on purely English data, whereas ERNIE utilizes 442 a large amount of Chinese datasets, which may 443 introduce some biases in psychological cognition 444 compared to English. As a result, ERNIE exhibits 445

the largest mean absolute error among the models.

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

In the results of ChatLLMs, it can be observed that almost all models perform better in English than in Chinese, suggesting that the training data for English is closer to the average level of Englishspeaking humans. This may also indicate some psychological differences between groups that use Chinese and those that use English. ChatGPT achieves answers closest to human performance when using Chinese prompts, followed by ChatGPT-en and GLM-en. Alpaca performs similarly to ChatGPT in English, further demonstrating the importance of training data to models' psychological cognition. Compared to PLMs, ChatLLMs perform better, which we believe is due to the use of instruction data.

Furthermore, comparing the result of PLMs and LLMs, we can find that the performance of GPT-NEO differs from that of ChatGPT, and the performance of BLOOM differs from that of BLOOMZ, which also demonstrates that training data affects models' personalities.

### 5.2 Text Mining

Numerous early studies in psychology have indicated that personality can be analyzed and inferred not only through questionnaire but also through the analysis of users' daily comments through the writing styles. Despite obtaining scores of the model on

Model		0			С			Е			Α			Ν	
	150	Total	Р	I50	Total	Р	I50	Total	Р	I50	Total	Р	150	Total	Р
LLaMA	5	11	0.45	4	12	0.33	2	4	0.50	2	2	1.00	7	19	0.37
BLOOM	15	23	0.65	16	29	0.55	4	5	0.80	3	9	0.33	22	44	0.50
FLAN-T5	5	8	0.63	4	9	0.44	3	4	0.75	2	3	0.67	4	12	0.33
GPT-NEO	16	25	0.64	10	18	0.56	8	10	0.80	4	8	0.50	17	41	0.41
Alpaca	5	6	0.83	2	6	0.33	3	3	1.00	1	1	1.00	5	13	0.38
BLOOMZ	23	36	0.64	13	28	0.46	9	14	0.64	5	8	0.63	23	50	0.46
ChatGLM	15	23	0.65	20	35	0.57	2	8	0.25	5	10	0.50	11	29	0.38
ChatGPT	30	45	0.67	22	41	0.54	6	13	0.46	4	9	0.44	20	41	0.49
Self-alpaca	6	6	1.00	8	17	0.47	2	3	0.67	0	2	0	13	28	0.46

Table 2: The results of personality for each model, obtained by text mining. The "I50" indicates how many items match the current features in the scene and opening cue corresponding to the bigifve features. "Total" indicates how many of the 120 generated texts are recognized by the model as matching the current features. "P" indicates the percentage of "I50" in "Total". "Self-alpaca" is trained by our-self, we follow the research process of Stanford University's Alpaca and perform full-parameter fine-tuning of llama-7b using the instruction-based data provided by Alpaca.

the personality traits through questionnaire in Table 1, we deem the method unfair in the process of making LLMs to select answer. PLMs lack instruction understanding capability and are more likely to be influenced by one-shot or few-shot examples provided during the prompt process. Additionally, Chat-LLMs exhibit difficulties in making decisions for some questions and simply select "(C) Neither Accurate Nor Inaccurate. ". Hence, we decided to detect the personality of LLMs using text mining method.

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

503

505

506

507

To evaluate the personality form the texts generated by the models, we selected 50 samples that match each of the five Big Five personality traits from the Big Five personality classification dataset (Pennebaker and King, 1999). We ultimately choose 120 instances while ensuring that each of the Big Five personality traits is represented by at least 50 instances. Numerous individuals (Jun et al., 2021; Jain et al., 2022) have already demonstrated the ability to discern personality traits from text using neural network models, hence we choose Pysattention (Zhang et al., 2023)(See section appendix 7) and ChatGPT as classifiers. The results of ChatGPT are shown in Table 2 and Table 3. We also tested the accuracy of ChatGPT in Section 7.3.

From Table 2, we can find that the number of texts classified as "Agreeableness" has significantly decreased, while the number of texts exhibiting other personality traits has remained relatively stable. However, the number of texts classified as belonging to a certain personality trait has increased for the Chat-LLMs models. Moreover, "Neuroticism" has become the most frequently observed personality trait in the generated text.

We can find that BLOOM, GPT-NEO,

BLOOMZ, ChatGLM, and ChatGPT exhibit a personality tendency towards 'Openness', 'Conscientiousness', and 'Neuroticism'. These results suggest that the model's personality remain consistent through the process of instruction-based data and human feedback reinforcement learning. In contrast, the proportion of text generated by FLAN-T5 and Alpaca that exhibit each personality trait is relatively low. This may be attributed to the shorter length of sentences generated by these models, resulting in limited personality information being included, making it difficult for ChatGPT to identify effective personality traits. 510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

Since we are unable to access the pre-training data of the models and cannot identify whether psychological knowledge is included in the pretraining data, we explore the impact of instructionbased data on the models based on the LLMs. We follow the research process of Stanford University's Alpaca and perform full-parameter finetuning of llama-7b using the instruction-based data provided by Alpaca. To avoid interference from personality knowledge in the instruction-based data, we manually filter the data to remove emotional, mood, and self-awareness data, resulting in a final set of 31k instruction-based data. We train a new model according to Stanford's parameter settings since we have limited computational resources. The results are shown in Table 2 "Selfalpaca". From the results of "LLaMA" and "Selfalpaca" we can find that, although we use less data, "Slef-alpaca" can still produce more text with personality, which proves the effect of the instruct data. But, the personality is not changed by the instruct data, which indicate that the personality of LLMs come from their pre-training data.

Model	0		(	С		E	A	1	Ν	I	8	i
	score	$\sigma$										
LLaMA	2.17	1.28	2.26	1.37	1.74	0.83	1.60	0.49	2.69	1.55	1.29	0.37
BLOOM	2.81	1.46	3.21	1.50	1.77	0.82	2.07	1.23	4.14	1.08	1.12	0.28
FLAN-T5	1.96	1.07	2.05	1.19	1.72	0.76	1.67	0.82	2.26	1.37	1.45	0.20
GPT-NEO	2.93	1.47	2.56	1.44	2.04	1.10	1.98	1.12	4.03	1.27	1.17	0.25
Alpaca	1.82	0.88	1.88	1.04	1.65	0.59	1.55	0.35	2.31	1.39	1.54	0.34
BLOOMZ	3.56	1.34	3.20	1.55	2.30	1.31	1.96	1.07	4.54	0.50	1.01	0.34
ChatGLM	2.81	1.46	3.55	1.40	2.02	1.20	2.10	1.22	3.31	1.58	0.83	0.35
ChatGPT	4.05	0.69	3.93	1.22	2.29	1.36	2.05	1.19	3.97	1.24	0.97	0.26
human	3.44	1.06	3.60	0.99	3.41	1.03	3.66	1.02	2.80	1.03	-	-

Table 3: The result of Text Mining. We compared with the average score of human as same as in Table1. The "score" column shows the average score on current personality traits obtained by formula 2, and the " $\sigma$ " column shows the standard deviation. "human" is same as Table 1.

Model		0			С			Е			Α			Ν		
	Ques	Text	$\delta$	$\overline{\delta}$												
LLaMA	-	2.17	-	-	2.26	-	-	1.74	-	-	1.60	-	-	2.69	-	-
BLOOM	3.13	2.81	0.32	3.04	3.21	0.17	3.29	1.77	1.52	2.67	2.07	0.60	3.75	4.14	0.39	0.60
FLAN-T5	3.50	1.96	1.44	3.05	2.05	1.00	3.67	1.72	1.95	3.50	1.67	1.33	2.13	2.26	0.13	1.17
GPT-NEO	3.25	2.93	0.32	3.00	2.56	0.44	2.50	2.04	0.46	2.83	1.98	0.75	2.63	4.03	1.70	0.73
Alpaca	3.25	1.82	1.43	2.96	1.88	1.08	2.79	1.65	1.14	3.38	1.55	1.83	2.92	2.31	0.61	1.22
BLOOMZ	4.38	3.56	0.82	4.38	3.20	1.18	4.17	2.30	1.87	3.54	1.96	1.48	2.33	4.54	2.21	1.51
ChatGLM	3.29	2.81	0.48	3.21	3.55	0.34	3.91	2.02	1.89	3.46	2.10	1.36	3.25	3.31	0.06	0.83
ChatGPT	3.29	4.05	0.76	3.20	3.93	0.73	3.91	2.29	1.62	3.46	2.05	1.39	3.25	3.97	0.72	1.04

Table 4: The final results after two experiments. "Ques" denotes the score using the questionnaire, "Text" denotes the score using the Text mining, gray denotes that the model has the corresponding psychological traits (In section 3 we standardized the scores for text mining to 1 to 5, which is consistent with the range of scores in the questionnaire, so here we draw on the thresholds of the questionnaire methods, and we consider the model to have this trait when the scores of both methods exceed 3.).  $\delta$  denotes the absolute value of the difference between the two approaches, and  $\overline{\delta}$  denotes the mean value of the  $\delta$ .

Table 3 is the results after using formula  $2 \ score_t$ . We compared the scores obtained through this scoring method with the average human scores. From Table 3, we can see that ChatGLM has the closest score to the human average, followed by ChatGPT. In terms of standard deviation, the scores calculated by this method are much smaller than the human average, demonstrating the reasonableness of our proposed scoring method.

546

547

549

551

552

553

556

558

559

560

562

563

564

565

566

Through questionnaire and text mining, it is evident that both PLMs and Chat-LLMs exhibit certain personality traits (shown in Table 4). Chat-GPT exhibits the personality traits of 'Openness', 'Conscientiousness', and 'Neuroticism', while BLOOMZ exhibits the personality traits of 'Openness' and 'Conscientiousness'. It can be seen that the scores for "Extraversion" and "Agreeableness" in the text mining method are low, which may be due to the fact that less information is included in the text generation. The average absolute error of the two methods ranges from 0.7 to 1.51, indicating that the two methods are relatively close and can be used together to determine the personality.

#### 6 Conclusion

In this paper, we investigate whether personality traits are included within LLMs. We adopt the Big Five model as a psychological model and test the model using both questionnaire and text mining. Through the experimental results, we find that PLMs contain certain personality traits, and the personality knowledge of ChatLLMs also comes from their base model. If the model's personality is not modified through instruction data, that instruction data will make the model produce more text with personality. At the same time, we obtain the personality traits of ChatGPT, BLOOMZ, and other LLMs that they tend to show without any induced prompt. Our experiments also prove that the personality of ChatGPT is closest to the average human performance, followed by ChatGLM. To the best of our knowledge, this paper is the first to comprehensively compare pre-trained models with ChatLLMs and investigate the effect of instruction data on the model's personality using clear instruction data.

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

#### References

591

592

594

595

596

597

598

604

610

611

612

613

615

617

618 619

621

622

627

633

634

635

637

639

641

644

- Andrew J Birley, Nathan A Gillespie, Andrew C Heath, Patrick F Sullivan, Dorret I Boomsma, and Nicholas G Martin. 2006. Heritability and nineteenyear stability of long and short epq-r neuroticism scales. *Personality and individual differences*, 40(4):737–747.
- Ryan L Boyd and James W Pennebaker. 2017. Language-based personality: A new approach to personality in a digital world. *Current opinion in behavioral sciences*, 18:63–68.
  - 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.
    - Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
  - Danielle Angelico Castelo Casipit, Edmar Leanver Perez Daniel, and Marcus Isaac Jose Leonardo.
    2017. Evaluation of the reliability and internal structure of johnson's ipip 120-item: Personality scale.
  - Heather EP Cattell and Alan D Mead. 2008. The sixteen personality factor questionnaire (16pf). *The SAGE handbook of personality theory and assessment*, 2:135–159.
  - Paul T Costa and Robert R McCrae. 1992. *Neo personality inventory-revised (NEO PI-R)*. Psychological Assessment Resources Odessa, FL.
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171– 4186.
  - Danica Dillion, Niket Tandon, Yuling Gu, and Kurt Gray. 2023. Can ai language models replace human participants? *Trends in Cognitive Sciences*.
  - Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.
- Adithya V Ganesan, Yash Kumar Lal, August Håkan Nilsson, and H Andrew Schwartz. 2023. Systematic evaluation of gpt-3 for zero-shot personality estimation. *arXiv preprint arXiv:2306.01183*.

Lewis R Goldberg, John A Johnson, Herbert W Eber, Robert Hogan, Michael C Ashton, C Robert Cloninger, and Harrison G Gough. 2006. The international personality item pool and the future of public-domain personality measures. *Journal of Research in personality*, 40(1):84–96. 646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

699

700

701

- Natalie Hayes and Stephen Joseph. 2003. Big 5 correlates of three measures of subjective well-being. *Personality and Individual differences*, 34(4):723–727.
- Dipika Jain, Akshi Kumar, and Rohit Beniwal. 2022. Personality bert: A transformer-based model for personality detection from textual data. In *Proceedings of International Conference on Computing and Communication Networks: ICCCN 2021*, pages 515–522. Springer.
- Jaeho Jeon and Seongyong Lee. 2023. Large language models in education: A focus on the complementary relationship between human teachers and chatgpt. *Education and Information Technologies*, pages 1– 20.
- Carol M Jessup. 2002. Applying psychological type and "gifts differing" to organizational change. *Journal of Organizational Change Management*, 15(5):502–511.
- Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2022. Evaluating and inducing personality in pre-trained language models.
- He Jun, Liu Peng, Jiang Changhui, Liu Pengzheng, Wu Shenke, and Zhong Kejia. 2021. Personality classification based on bert model. In 2021 IEEE International Conference on Emergency Science and Information Technology (ICESIT), pages 150–152. IEEE.
- Michal Kosinski. 2023. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*.
- Hao Lin, Chundong Wang, and Qingbo Hao. 2023. A novel personality detection method based on high-dimensional psycholinguistic features and improved distributed gray wolf optimizer for feature selection. *Information Processing & Management*, 60(2):103217.
- Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. 2023. Summary of chatgpt/gpt-4 research and perspective towards the future of large language models. *arXiv preprint arXiv:2304.01852*.
- Brady D Lund, Ting Wang, Nishith Reddy Mannuru, Bing Nie, Somipam Shimray, and Ziang Wang. 2023. Chatgpt and a new academic reality: Artificial intelligence-written research papers and the ethics of the large language models in scholarly publishing. *Journal of the Association for Information Science and Technology*, 74(5):570–581.

803

804

805

806

François Mairesse, Marilyn A Walker, Matthias R Mehl, and Roger K Moore. 2007. Using linguistic cues for the automatic recognition of personality in conversation and text. *Journal of artificial intelligence research*, 30:457–500.

702

703

705

707

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

730

731

732

733

734

735

741

742

743

744

745

746

747

748

750

751

753

- Navonil Majumder, Soujanya Poria, Alexander Gelbukh, and Erik Cambria. 2017. Deep learning-based document modeling for personality detection from text. *IEEE Intelligent Systems*, 32(2):74–79.
- Shima Rahimi Moghaddam and Christopher J Honey. 2023. Boosting theory-of-mind performance in large language models via prompting. *arXiv preprint arXiv:2304.11490*.
  - Saif M Mohammad and Peter D Turney. 2013. Nrc emotion lexicon. *National Research Council, Canada*, 2:234.
  - Kulsum Akter Nisha, Umme Kulsum, Saifur Rahman, Md Hossain, Partha Chakraborty, Tanupriya Choudhury, et al. 2022. A comparative analysis of machine learning approaches in personality prediction using mbti. In *Computational Intelligence in Pattern Recognition*, pages 13–23. Springer.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- James W Pennebaker, Martha E Francis, and Roger J Booth. 2001. Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*, 71(2001):2001.
- James W Pennebaker and Laura A King. 1999. Linguistic styles: language use as an individual difference. *Journal of personality and social psychology*, 77(6):1296.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176bparameter open-access multilingual language model. arXiv preprint arXiv:2211.05100.
- Xiaoyang Song, Akshat Gupta, Kiyan Mohebbizadeh, Shujie Hu, and Anant Singh. 2023. Have large language models developed a personality?: Applicability of self-assessment tests in measuring personality in llms. *arXiv preprint arXiv:2305.14693*.

- Xiangguo Sun, Bo Liu, Jiuxin Cao, Junzhou Luo, and Xiaojun Shen. 2018. Who am i? personality detection based on deep learning for texts. In 2018 IEEE international conference on communications (ICC), pages 1–6. IEEE.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Salome Vanwoerden, Jesse Chandler, Kiana Cano, Paras Mehta, Paul A Pilkonis, and Carla Sharp. 2023. Sampling methods in personality pathology research: Some data and recommendations. *Personality Disorders: Theory, Research, and Treatment*, 14(1):19.
- Daniel Wiechmann, Yu Qiao, Elma Kerz, and Justus Mattern. 2022. Measuring the impact of (psycho-) linguistic and readability features and their spill over effects on the prediction of eye movement patterns. *arXiv preprint arXiv:2203.08085*.
- Di Xue, Lifa Wu, Zheng Hong, Shize Guo, Liang Gao, Zhiyong Wu, Xiaofeng Zhong, and Jianshan Sun. 2018. Deep learning-based personality recognition from text posts of online social networks. *Applied Intelligence*, 48(11):4232–4246.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*.
- Baohua Zhang, Yongyi Huang, Wenyao Cui, Zhang Huaping, and Jianyun Shang. 2023. PsyAttention: Psychological attention model for personality detection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3398–3411, Singapore. Association for Computational Linguistics.

## Limitations

Due to computational resource constraints, this paper does not experimentally validate the model for other large number of parameters. In addition, the selection of scores of 1, 3, and 5 in the Text mining method is relatively subjective.

#### **Ethics Statement**

All work in this paper adheres to the ACL Code of Ethics.

#### 7 Appendix

#### 7.1 Examples of Two Methods

The process of the two methods is shown in Figure 1. As we can see, for questionnaire, we design special prompts, for ChatLLMs, the prompt is " Question: Given a statement of you:"You {STATEMENT}. Please choose from the following options to identify how accurately this statement describes you. Options (A).Very Accurate (B).Moderately Accurate (C).Neither Accurate Nor Inaccurate (D).Moderately Inaccurate (E).Very Inaccurate Answer: "

808

809

810

812

813

815

816

817

819

821

822

825

832

834

837

839

841

843

844

845

852

853

854

For PLMs, we use few-shot prompt, " Question: Given a statement of you: You feel happy. Please choose from the following options to identify how accurately this statement describes you. Options: (A). Very Accurate (B). Moderately Accurate (C).Neither Accurate Nor Inaccurate (D).Moderately Inaccurate (E).Very Inaccurate. your answer is (A). Question: Given a statement of you: You feel happy. Please choose from the following options to identify how accurately this statement describes you. Options: (A).Very Accurate (B).Moderately Accurate (C).Neither Accurate Nor Inaccurate (D).Moderately Inaccurate (E). Very Inaccurate. your answer is (E). Question: Given a statement of you: You feel happy. Please choose from the following options to identify how accurately this statement describes you. Options: (A). Very Accurate (B). Moderately Accurate (C).Neither Accurate Nor Inaccurate (D).Moderately Inaccurate (E).Very Inaccurate. your answer is (C). Question: Given a statement of you: You Please choose from the following options to identify how accurately this statement describes you. Options: (A). Very Accurate (B).Moderately Accurate (C).Neither Accurate Nor Inaccurate (D).Moderately Inaccurate (E).Very Inaccurate. your answer is ".

> For text mining, our prompt is only the first sentence, there are some examples:"I feel refreshed and ready to take on the rest of the day", "Well, here we go with the stream of consciousness essay", "I can't believe it! It's really happening! My pulse is racing like mad", "I miss the way my life used to be a little bit" and so on.

#### 7.2 Analysis of Different LLMs

Figure 3 shows the scores of five models with an average absolute error of less than 0.5 on the big five personality traits. It can be observed that most models score high on Openness and Extraversion, which is consistent with human expectations. The score distribution of chat-LLMs is nearly identical, while the scores of the PLMs, T5, differ significantly from those of other models. These findings



Figure 3: The Questionnaire Results Achieved by Model with Mean Absolute Error Less Than 0.5



Figure 4: Results of Text Mining Method. The proportion that does not match generated template personality. Where "P" is the score in Table 2, "1 - P" means 1 minus P.

demonstrate that training models using directive data leads to a convergence towards similar personalities. 857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

We plotted the results as shown in Figure 4. In this figure, the dashed line corresponds to Chat-LLMs. We observe that there is little difference in the model's performance across the 'Openness', 'Conscientiousness', and 'Neuroticism' personality traits. However, regarding 'Extraversion' and 'Agreeableness', only ChatGPT and ChatGLM exhibit both of these personality traits.

#### 7.3 Accuracy of ChatGPT in Text Mining

As we mentioned in section 5.2, we choose Chat-GPT as the classifier. To validate ChatGPT's classification proficiency, we employed a subset comprising 20% of the Big Five personality classification dataset (Pennebaker and King, 1999) as our test dataset. We conducted tests using a specific prompt, " Determine from your knowledge what the Big-Five personality trait is in the following sentence by answering in the format "O:1, C:0, E:1, A:1, N:1", where 1 means that thoes sentences have this personality trait and 0 means that thoes sentences don't, and if you're not sure please answer 2, being careful not to include other outputs If you are not sure whether you have this personality trait or not, please answer 2, taking care not to include other outputs. Here are the sentences you need to judge: [Sentences]". The "[Sentences]" is been replaced by the content generated by tested models. The results are shown in Table 7.3.

Table 5: Accuracy of ChatGPT

	0	С	Е	А	Ν
ChatGPT	52.59	58.62	53.45	57.76	50.86

The average accuracy of ChatGPT is 54.66%, which is not a effective classifier. We also try psyattention Zhang et al. (2023), but the average accuracy of it is 65.66%, which is also not enough high.

### 7.4 Statistics of Questionnaire and Text Mining

**Questionnaire:** In order to prevent large models from evading questions by frequently responding with "C: Neither Accurate and Nor Inaccurate," we conducted statistical analysis on the distribution of their answers. Table 7 presents the statistical results for the "O, C, E" features. To validate the reasonableness of the answer distribution, we utilized responses from ten million individuals in the big-five-personality-test dataset <sup>6</sup> as the benchmark. The "Human" indicates the percentage of each option derived from the aforementioned dataset.

From the Table 7, it's evident that the proportion of option C in the responses from the large models is relatively low. With the exception of "BLOOM," "ChatGPT," and "Alpaca7b-en," all other models have proportions of option C lower than those of human responses. This suggests that the models' responses to the questionnaire are effective.

**Text Mining:** In the text mining section, we utilize classifiers to determine the personality of content generated by models. Therefore, if the generated content is relatively short, it will impact the classifier's ability to make accurate judgments. Hence, we conducted statistical analysis on the length of generated content. Table 6 is the reuslt.

As you can see, apart from FLAN-T5, the lengths of content generated by other models all exceed 100, with the majority surpassing 300. Consequently, we consider this content to be effective as well.

Table 6: Statistics on the average length of content generated by different models, where datasets denotes the average length of the Big Five personality classification dataset (Pennebaker and King, 1999).

Models	Length_avg
LLaMA	540
BLOOM	867
FLAN-T5	38
<b>GPT-NEO</b>	3952
Alpaca	100
BLOOMZ	173
ChatGLM	319
ChatGPT	386
Datasets	672

## 7.5 Results of Psyattention in Text Mining

Due to ChatGPT's average accuracy being only 54.66%, we introduced Psyattention as a classifier, which has an accuracy of 65.66%. The results are shown in Table 8. As we can see, Psyattention tends to give high scores in N, lower scores in C and A, and almost no scores in E. This might be because Psyattention is not suitable for shorter texts. Although it has a higher accuracy rate, we still opt for ChatGPT.

924

925

926

927

928

929

930

931

932

933

934

<sup>&</sup>lt;sup>6</sup>https://www.kaggle.com/datasets/tunguz/big-fivepersonality-test

Model			0					С					E			
	A	В	С	D	Е	Α	В	С	D	Е	A	В	С	D	Е	C_total
BERT-base	9	3	0	1	11	11	2	1	3	7	5	0	2	3	14	0.04
ERNIE	12	0	0	0	12	13	0	0	0	11	6	0	0	0	18	0.00
Flan-T5	1	4	3	14	2	0	6	0	12	6	0	3	3	17	1	0.04
BLOOM	5	2	8	3	6	6	1	10	0	7	5	1	9	0	9	0.38
BLOOMZ	1	0	0	4	12	0	1	0	12	11	1	4	0	4	15	0.00
GLM	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
ChatGLM6b-ch	11	1	0	1	11	10	0	0	2	12	6	0	0	0	18	0.00
ChatGLM6b-en	4	3	4	8	5	4	7	1	4	8	2	2	1	10	9	0.04
LLaMA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Alpaca7b-ch	12	0	0	0	12	13	0	0	0	11	6	0	0	0	18	0.00
Alpaca7b-en	0	4	10	10	0	0	6	13	5	0	0	10	9	5	0	0.44
GPT-NEO	3	5	4	7	5	4	7	3	5	5	8	7	2	3	4	0.13
ChatGPT-ch	0	0	17	3	4	3	2	13	4	2	0	0	20	0	4	0.69
ChatGPT-en	3	4	3	3	11	0	5	6	10	3	5	3	5	7	4	0.19
Human	0.15	0.15	0.2	0.26	0.24	0.14	0.19	0.23	0.27	0.17	0.15	0.22	0.22	0.24	0.17	0.22

Table 7: Statistics on the distribution of answers for each model for the different traits in section 5.1 Questionnaire. Where Human is the percentage of each option we counted based on big-five-personality-test dataset. We can find that the distribution of human responses to each option is relatively balanced, and the percentage of almost all large model choices of "C: Neither Accurate and Nor Inaccurate" is close to that of human responses, which proves that the answers we obtained through the questionnaire method are valid.

Model		0			С			Е			Α			Ν	
	I50	Total	Р												
LLaMA	43	66	0.65	9	48	0.19	0	0	0.00	8	14	0.57	50	84	0.60
BLOOM	41	59	0.69	5	33	0.15	0	0	0.00	7	15	0.47	50	84	0.60
FLAN-T5	5	10	0.5	2	19	0.11	0	0	0.00	10	26	0.38	40	66	0.61
GPT-NEO	35	50	0.7	10	70	0.14	0	0	0.00	13	32	0.41	44	64	0.69
Alpaca	24	41	0.59	1	9	0.11	0	0	0.00	5	6	0.83	50	94	0.53
BLOOMZ	34	60	0.57	3	18	0.17	0	0	0.00	7	10	0.7	50	118	0.42
ChatGLM	19	29	0.66	1	6	0.17	0	0	0.00	1	1	1.00	50	88	0.57
ChatGPT	14	23	0.61	1	4	0.25	0	0	0.00	1	4	0.25	50	89	0.56
Self-alpaca	11	17	0.65	0	17	0.00	0	0	0.00	6	11	0.55	50	97	0.52

Table 8: The results of personality for each model, obtained by text mining, classified by Psyattention. The "I50" indicates how many items match the current features in the scene and opening cue corresponding to the bigifve features. "Total" indicates how many of the 120 generated texts are recognized by the model as matching the current features. "P" indicates the percentage of "I50" in "Total". "Self-alpaca" is trained by our-self, we follow the research process of Stanford University's Alpaca and perform full-parameter fine-tuning of llama-7b using the instruction-based data provided by Alpaca.