PersonaMovs: A Multimedia Conversational Personality Dataset

Anonymous ACL submission

Abstract

Automatic personality detection has evolved from simple text classification to sophisticated multimodal analysis, recognizing the multidimensional manifestation of personality beyond textual data. This shift highlights the need for datasets that can accurately capture the complexity of human personality through diverse modalities. We introduce the PersonaMovs, a large, extensive and varied multimedia conversational dataset, built on 305 movies and 14 TV series, featuring over 46k dialogues, 552k utterances, 4016 characters, and 963 hours of video. PersonaMovs not only addresses the challenges of existing datasets by offering majority-voted personality annotations and detailed social relation networks but also paves the way for advanced analysis of interactions of personality across various contexts.

1 Introduction

001

006

007

011

012

017

019

024

027

Personalities is a comprehensive yet complex trait that encapsulates individual differences in patterns of thinking, feeling, and behaving (Costa and Mc-Crae, 2002). Detecting personality automatically is of significant importance for improvement of machine's ability to have human-like cognition and engage in more natural interactions with humans, particularly in the context of advancing Artificial General Intelligence (AGI) and various practical applications such as reflective linguistic programming (Fischer, 2023), disease diagnosis (Tseng et al., 2013) and mental health prediction (Feng et al., 2024). At the very beginning, owing to the limitations of multimedia model and computational power, researchers initially approached personality prediction as a straightforward text classification task, focusing on deciphering personality traits from individuals' digital footprints online (Kerz et al., 2022; Yang et al.). However, as illustrated in Figure 1, there is now

growing recognition that personality is expressed across multiple modalities, with nuances that cannot be fully captured through text-based analysis alone (Al Maruf et al., 2022; Zhu et al., 2022; Bose et al., 2023). This realization has driven the shift towards multimodal personality detection as the dominant methodology.



Figure 1: The Distinctive Features in Three Modalities for Personality Prediction

Personality datasets, integrating text, audio, visual information along with the manner of speaking, face expressions, body language and so on, offer a richer, more nuanced view of human behavior and personality expressions than text-based datasets alone. This comprehensive approach is vital for developing models that accurately cap-

051

ture the complexity of human personality. Consequently, numerous multimodal datasets have been released in recent years, and several efforts have focused on constructing such datasets specifically for personality analysis (Palmero et al., 2021; Junior et al., 2021; Jiang et al., 2020; Chen et al., 2022). Additionally, many multimodal datasets designed for other tasks have been adapted for personality prediction. For example, TVQA (Lei et al., 2018), a large-scale dataset originally created for visual question answering, is frequently utilized in personality research due to its extensive scope.

057

061

062

067

077

079

081

094

100

101

102

103

104

105

Recently Li et al.; Pal et al. state that personality is far beyond discrete and reveal that dynamic nature of personal identity. Although current datasets have evolved to include many features necessary for personality prediction, they still exhibit several limitations:

1) **Limited Data Source**: Previous datasets often select one or several famous movies or TV series as the raw data, resulting in a limited number of characters and personality types covered, which hinders the generalizability of model training.

2) **Manual Annotation Issues**: The process of manually annotating personality traits typically relies on a few numbers of volunteers with varying levels of expertise, leading to potential inconsistencies and biases in the annotations.

3) Lack of Relations Representation: Current datasets largely fail to account for the interactions of personality, which evolves over time in response to changes in life experiences and social environments. These datasets often provide static snapshots of personality traits, overlooking the psychological understanding that personality can vary across different contexts and stages of life.

In our study, we endeavor to partially eliminate the aforementioned limitations by providing a scale-up multimodal dataset that contains reliable labels. Specifically, we find a personality database website¹ that offers a large amount of personality types for virtual characters and (Zhu et al., 2023) have scraped the personality data from it to annotate TVQA dataset. Compared with previous datasets whose labelling commonly involved five to ten people, our datasets are labelled by about 160 voters on average. It shows the vote distribution rather than a single personality type which is more persuasive and operable.

Against these backdrops, we introduce the PersonaMovs, a comprehensive dataset that starkly contrasts with existing offerings in several key aspects. PersonaMovs is built on 305 movies and 14 TV series (894 episodes in total) in different genres, including more than **46k dialogues**, **552k utterances**, **4016 characters** and **963 hours videos**. With the rich annotation, our dataset supports 4 personality traits models (MBTI, Big Five, Enneagram and Instinctual Variant) with more 28 personality classes, 7 kinds of Social Relations. Our analysis highlights substantial quantity and diversity in content, adequate experiments on different models with all modalities and personality interactions discovery. 106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

Our contributions are as follows:

- We introduce PersonaMovs, the most comprehensive and varied multimodal personality dataset to date, surpassing existing datasets in **scope** and **diversity**. This dataset uniquely combines movie and TV genres with personality analysis via audio, video, and text, along with crowd sourced personality and social relations labels, unlocking new avenues in personality research².
- We study seven model architectures from different model families. Our results show that PersonaMovs is more **difficult** compared to other datasets, not only because it has a larger amount of multimedia data, but also due to its diversity and similarity to real life. Through our analysis, we found that our dataset provides a more realistic reflection of the distribution of personalities in the real world, further enhancing its value for personalityrelated tasks.
- For the first time, we categorize 7 types of social relations to depict character interactions on a scene-by-scene basis, enabling a granular analysis of personality through relations networks. Guided by the relations networks, we identify psychological phenomenons in conversational contexts, which largely explain the interaction of personality statistically.

¹https://www.personality-database.com/

²https://anonymous.4open.science/r/ sample-of-MMPD-0177

Dataset	Dialogues	Utters. / Dial	Characters	Source
MEmoR	8.53k	64.23	7	The Big Bang Theory
FriendsPersona	0.71k	27.61	7	Friends
CPED	12k	1	392	40 TV shows
UDVIA	188	65.31	147	Dyadic Interaction
The ChaLearn FI	10k	Unknown	3000	YouTube
TVQA	29.4k	2.2	Unknown	6 TV shows
PersonaMovs	46.21k	12.42	4000+	300+ Movies and 14 TV Shows

Table 1: Comparison of different datasets and our PersonaMovs

2 Dataset Design

152

154

155

156

157

158

159

160

162

163

164

165

166

167

168

169

170

2.1 Personality Theory

Personality refers to the combination of characteristics or qualities that form an individual's distinctive character. It encompasses a wide range of traits, behaviors, thoughts, and emotional patterns that evolve from biological and environmental factors (Lepri et al., 2012). A particular personality can determine various outward observable properties or features, including consistent behavioral patterns, communication style, emotional expression and so on. These traits manifest in how an individual consistently acts and reacts in different situations, their manner of speaking and body language, the openness or restraint of their emotional displays, their ways of relating to others, their approach to making decisions, and their preferences in activities, hobbies, and social engagements.

2.2 Multiple Personality Models are Needed

In constructing such a dataset for personality 171 prediction, incorporating four distinct personal-172 ity models, provides a comprehensive framework 173 for understanding the multifaceted nature of hu-174 man personality. To this end, evolving four distinct 175 personality models-Myers-Briggs Type Indica-176 tor (MBTI), Big Five, Enneagram, and Instinc-177 tual Variant-into our dataset construction is es-178 sential. Each of these models provides a unique 179 lens through which to view and interpret personal-180 ity traits, offering complementary insights that are critical for a holistic understanding(See all details 182 about each trait in Appendix A). For instances, MBTI has 16 unique personality types. Among these, INFJ is a personality type with the Intro-186 verted, Intuitive, Feeling, and Judging traits. They tend to approach life with deep thoughtfulness and 187 imagination. Their inner vision, personal values, and a quiet, principled version of humanism guide them in all things. 190

2.3 Relations to Interpret Personality

Human social networks are complex and multifaceted. By categorizing relations, we can better understand the nuances of how people interact with each other. Different types of relations provide context for interactions, which is crucial for analyzing social behaviors and patterns, improving social network analysis, and applying this knowledge across various fields and applications. 191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

225

226

227

229

230

Each personality model offers unique insights and covers different aspects of personality, making them collectively valuable for a multidimensional approach to personality prediction. In addition to these models, we introduce 7 categories of social relations among characters. These provide a comprehensive framework to observe and interpret the nuances of personality in action. We identify seven types of social relations from the perspectives of psychology and sociology. Based on various social environments, we categorize 7 kinds of social relations (shown in Table 2) including *family*, *friendship*, *romantic*, *professional*, *social*, *academic and online*.

2.4 Structure of PersonaMovs

Aiming to deliver a tidy and readable structure, we distribute different scenes in a single text file with corresponding audio and video clip. As shown in Figure 2, it provides a timestamp of the scene, a visual snapshot from the movie, and a transcription of their dialogue. The characters' personalities are profiled using various typologies, such as MBTI and Enneagram, with accompanying bar charts showing the distribution of votes or consensus on these personality assessments. Additionally, the social relations between Gatsby and Daisy is categorized as romantic, offering insights into their interaction.

3 Methodology

In this section, we outline the methodology employed to gather, process, and annotate the data

Relations type	Description
Family	Parents (grandparents) and children, siblings, etc.
Friendship	Based on common interest, mutual respect and affection, but not related to the blood.
Romantic	Based on emotional attraction and include dating, marriage, etc.
Professional	Formed in a work environment, such as colleagues, superiors and subordinates, etc.
Social	Formed in a broader social context, such as neighbors, club members.
Academic	Formed in an educational setting, such as between teachers and students, classmates.
Online	Established in online spaces or through social media platforms.

Table 2: Descriptions of social relations



Figure 2: A sample from PersonaMovs, which includes full modality data along with timestamps, personality label, relationship label, and personality vote distribution.

for our study. We begin by detailing the sources of our multimodal video data and personality labels, focusing on how we efficiently align subtitles with original scripts to ensure accurate temporal and character associations. And we also present our annotation process, explaining how we leverage the ChatGPT to automatically annotate social relations among characters within the text data.

3.1 Source of Data

232

240

241

243

244

245

246

247

Our data source contains mainly two parts, the multimodal video data and personality labels. For video data, we include 14 different genres of TV series and movies via an open-source website¹, and for the scripts and subtitles, we also find other open-source websites²³ for research offering the free scripts and subtitles of many famous movie and television programs. Considering the insuffi-

cient labeling method of existing works, we collect the personality annotations from personality database website as well as the voting distribution and align them to correctly scripts.

248

249

252

253

254

255

256

257

258

261

262

263

264

265

266

3.2 Data Alignment Process

As subtitle contain temporal information and original scripts associate utterances with characters, we are supposed to align them properly as efficient as possible. However, most of the existing multimodal datasets annotate the timestamps manually with taking up a great deal of time. There are also some works which utilize different automatic tools to align the utterances with their corresponding information. For instance, (Lian et al., 2024) use an Automatic Sound Recognition (ASR) tool called Gentle⁴ to get the timestamps for the utterances. To streamline the process of aligning dialogue utterances with their respective timestamps and speakers from subtitles, we propose an effi-

¹https://yts.mx/

²https://www.simplyscripts.com/

³https://subscene.com/

⁴https://github.com/lowerquality/gentle

cient method leveraging a fuzzy matching algorithm (details in Appendix B). Following successful alignment, we proceed to segment the video content into distinct scenes according to the timestamps. Besides, we use $FFmpeg^1$ to extract the audio track from the video clips and output it as a .mp3 file.

267

268

272

273

274

277

278

279

284

288

294

295



Figure 3: Process of data alignment

3.3 Annotation Process

In this study, we propose a method to automatically annotate social relations among characters in scripted text using the GPT-3.5 model (OpenAI, 2023), which is well-suited for processing natural language data. Our approach begins by preprocessing the text and dividing it into scenes. For each scene, we design a specific prompt (Figure 4) that guide ChatGPT in identifying the social relations among characters.

The rationale for using ChatGPT in this context is based on the model's general understanding of the show, its characters, and the roles they play. Given that GPT has been trained on a vast corpus of text, it is able to infer the dynamics of social interactions, even without explicit annotations. This capability allows the model to effectively detect and label social relations in scenes, making it a plausible tool for this task.

4 Evaluation

We present the basic statistics of our dataset in the first part, and then we evaluate the accuracy



Figure 4: Prompt design for relations annotation

of our alignment and annotation process to ensure their reliability. Additionally, we not only test our dataset on different advanced models but also do ablation experiment on both modality and social relations annotation. To discover interesting topics about interaction of personality, we focus on those changes of personality and discover several interesting psychological phenomenons.

4.1 Dataset Statistics

As we mentioned before, PersonaMovs is not only a large dataset containing a huge amount of text, audio and video corpus but also its data is highly diverse in terms of personality types, movie and television production genres, and relationship types. Fig 5 are the distribution of MBTI types and social relations, which indicates the diversity in terms of diversity of personality labels and interaction scenarios.



Figure 5: Distribution of MBTI types and social relations

Algorithm Evaluation

To evaluate the performance of our character-tosubtitle matching algorithm, we manually check

296

297

298

299

301

302

303

314

315

¹https://ffmpeg.org/

the aligned characters' name based on the script.
The annotators were a group of five human volunteers with backgrounds in filmography and literary. They are in their mid-twenties and had at least
an undergraduate education. The protocol involves
the following steps:

324

325

327

328

331

332

334

338

359

- 1) Annotators independently reviewed a randomly selected sample of 50 dialogues from the dataset. By given relevant scripts and subtitle files, they are required to match each utterance in subtitle with corresponding names.
- 2) The annotations were compared against the results generated by our algorithm to evaluate accuracy.

The algorithm demonstrates an accuracy of about 88%, indicating a high level of accuracy in correctly identifying character names within subtitles across diverse content types. Compared to existing ASR matching algorithm, our approach gains an improvement by 5% in accuracy. Besides, our algorithm shows a very strong efficiency comparing the ASR method, of which accelerating almost 7 times.

Method	Movies	TV	Exec. Time (s)
Gentle (ASR)	82.71%	85.21%	26.51
Our algorithm	87.53%	88.98%	3.55

Table 3: Accuracy and running time per dialogue of subtitle matching algorithm

Annotation Accuracy

Despite ChatGPT's impressive capabilities, the la-341 342 bels annotated by it are not completely correct and its automated annotations require validation through human expertise. To measure the auto-345 matic annotation accuracy, we sampled 235 scenes randomly and involved 5 human labelers on relations annotation. These labelers are in their mid-347 twenties, undergraduate or higher education background, proficient in English with majors in psychology, filmography and sociology, who were instructed to select one of the designated social relations after aligned video. We continue to compare the automatically annotated results to the human-labeled ground truth. The outcome shows that both social relationship annotations are dependable, with the accuracy reaching aroud 98% and 93%.

> The dataset's foundation on crowdsourced voting allows for an in-depth analysis of subjec-

Relations type	Accuracy in Movies	Accuracy in TV
Family	99.18%	93.27%
Friendship	95.35%	92.65%
Romantic	95.21%	91.93%
Professional	98.29%	95.10%
Social	97.13%	94.07%
Academic	98.50%	91.54%
Online	99.00%	93.15%
Overall	98.21%	93.91%

Table 4: Accuracy of each category of social relations annotation

tive biases in personality perception. Researchers can investigate how different demographics (age, gender, cultural background) perceive personality traits and emotions in characters, revealing biases that may exist in personality assessment. This could also extend to studying the impact of viewer's own personality traits on their perceptions of characters, thus contributing to a deeper understanding of projection and identification processes in media consumption. 360

361

362

363

364

365

366

367

368

370

371

372

373

374

375

376

377

378

380

381

382

383

384

386

387

4.2 Experiment Results

Dataset Difficulties

We test our dataset on popular models including BERT (Devlin et al., 2019), D-DGCN (Yang et al., 2023), Roberta (Liu et al., 2019), AttR-CNN (Xue et al., 2018), GPT-3.5 (OpenAI, 2023), GPT-4 (OpenAI, 2024) and MCT (Sun and Zhang, 2023). We use the MBTI framework, which includes 16 distinct personality types, as the labels for our dataset. As shown in Table 5, the accuracy of our dataset is lower than that of other comparable datasets. Accuracy here refers to the proportion of correctly predicted personality types out of the total number of predictions. One of the key challenges contributing to this lower accuracy is the increased complexity and diversity of our dataset, which encompasses a broader range of multimedia content compared to other datasets.

Method	Modalities	FP	TVQA	PM
BERT	T only	61.14	60.61	52.94
D-DGCN	T only	69.56	70.21	68.47
Roberta	T only	62.58	69.24	60.37
AttRCNN	T only	65.01	67.25	62.44
GPT-3.5	T only	69.21	66.89	64.08
GPT-4	T & V	79.14	78.33	76.90
MCT	T, A & V	71.67	69.93	68.47

Table 5: Accuracy of different methods on Friends Persona (FP), TVQA, and PersonaMovs (PM). T, A, & V stand for text, audio and video respectively. Lowest accuracy in each row is bolded.

A more challenging dataset, such as the one we have developed, offers several advantages in 389 terms of personality detection: 1) Our dataset cap-390 tures a wide range of real-life situations and intricate contexts, which better mirrors the complexity of human interactions. This realism is crucial for developing models that can perform well in 394 practical applications. 2) Training on a more difficult dataset forces models to learn more nuanced patterns and relationships, leading to better generalization capabilities. 3) A difficult dataset sets a high standard for model evaluation, ensuring that only the most effective models are considered 400 successful. This helps in distinguishing truly ad-401 vanced models from those that perform well only 402 on simpler tasks. Additionally, one notable obser-403 vation from the results is that the MCT model, 404 which leverages three modalities (text, audio, and 405 video), does not outperform the GPT-4 model, 406 which uses only two modalities (text and video). 407 This performance gap suggests that Large Lan-408 guage Model outperforms the small model on this 409 task, even though the latter uses more modalities. 410

The Importance of Multi-Modality

411

419

420

421

422

423

424

425

426

427

428

429

412We conducted a series of ablation experiments to413assess the impact of different modalities and rela-414tions annotations on the performance of personal-415ity prediction models. The experiments were de-416signed to understand how the exclusion of spe-417cific modalities or relations annotations affects the418overall prediction accuracy.

Method	Modality	Accuracy
	T & A	66.13
MCT	T & V	67.91
MCI	T only	63.43
	T, A & V	68.47
CDT 4	T only	70.20
UF 1-4	T & V	76.90

Table 6: Ablation experiment on different modalities

Table 6 presents the results of ablation experiments where different combinations of video and audio modalities were excluded. The result underscore the critical importance of using multiple modalities to achieve higher accuracy in personality prediction tasks. Models that leverage both audio and video data, in addition to text, consistently outperform those that rely solely on textual data.

Table 7 shows the results of ablation experiments focusing on the inclusion or exclusion of social relation annotations which finds the relations

Method	With Relations	Without Relations
BERT	54.16	52.94
Roberta	59.04	58.39
GPT-4	74.49	71.20

Table 7: Ablation experiment on social relation annotations.

annotations tend to slightly enhance the performance. This highlights the importance of including rich contextual information to improve the accuracy of personality prediction models. 430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

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

The multimodal nature of the dataset enables comprehensive studies that integrate different data types to understand personality. This could lead to the development of new theories or the refinement of existing ones that account for the complexity of personality as depicted through various media. It could also foster interdisciplinary research, combining insights from psychology, computer science, linguistics, and media studies.

Personality Analysis

As mentioned in the Introduction, personality can change depending on various contexts. To explore this, we studied the relationship between the diversity of personalities and prediction accuracy. For each character, we calculate the entropy of the vote distribution to represent the uncertainty of personality as well as predicting their personalities using GPT-3.5.

Entropy, denoted as H(X), is a measure of the uncertainty or complexity in a probability distribution. It is calculated using the formula:

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_b(p(x_i))$$

where X is a random variable representing the personality distribution, $p(x_i)$ is the probability of each personality type x_i , and b is the base of the logarithm, typically 2 when measuring entropy in bits.

Figure 6 shows that as the complexity of a character's personality increases, the corresponding prediction accuracy decreases. This finding suggests that personality should be considered an attribute based on different context rather than a static one.

According to our finding, we conduct statistical analysis based on our dataset to figure out if there exists certain personalities that are easily attracted to each other. To analyze the patterns of personality attraction, we focus on identifying pairs of



Figure 6: Relationship between MBTI entropy and prediction accuracy

personalities that frequently appear together in ro-472 473 mantic and friendship relations. Figure 7 presents the favorite network with 16 MBTI personality 474 types, providing a clear visual summary of sta-475 tistical findings. The size of each node is propor-476 tional to the number of connections (degree) it has, 477 which means personality types with more relation-478 ships are represented by larger nodes. The color 479 of the edges represents the weight of the relation-480 ship between the personality types. Darker edges 481 indicate a higher frequency or stronger relation-482 ship. For instance, ESFP is the most popular per-483 sonality since almost every other personality has 484 a friendship relation with it, and ESTP prefers to 485 be around INFP, ESFP and people with the same 486 personality as themselves. Another interesting pat-487 tern emerges in ESTP relationships: while they 488 frequently form friendships with ESFPs, roman-489 tic relationships between these types are uncom-490 monly rare in our observations. Thus, these poten-491 tial findings may inspire the research on psychol-492 ogy and other social science. 493



Figure 7: Favorite betworks with different personalities about two relations

By including data on social relations between characters, the dataset opens new pathways for exploring the interaction of personality. This as-

494

495

496

pect provides a basis for computational models that simulate personality in social networks, potentially informing theories on social behavior, conflict resolution, and group dynamics. 497

498

499

500

501

502

503

504

505

506

507

508

509

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

5 Copyright Considerations

The movies and TV series incorporated in this dataset are under the copyright of their respective holders. They are utilized in this academic and research work in accordance with the fair - use principles as defined in relevant international copyright agreements such as the Berne Convention for the Protection of Literary and Artistic Works (Berne Convention) (Ber) and the Universal Copyright Convention (Uni), or based on specific permissions obtained from the copyright holders. It must be clearly stated that this usage does not imply any form of endorsement or affiliation with the copyright owners. The use of these materials is strictly restricted to the scope of the permission granted, and any commercial distribution is explicitly prohibited.

6 Conclusion

In this study, we introduce PersonaMovs, an outstanding multimodal dataset tailored for personality prediction. Built upon a foundation of varied movies and TV shows, PersonaMovs enriches with precise annotations for personality traits based on different psychological personality models and detailed relations networks, capturing the interplay of characters' interactions. By integrating multimodal data and emphasizing the nature of personality within social contexts, PersonaMovs opens new avenues for comprehensive analysis of personality interaction, offering valuable insights into how personality traits manifest and interact in varied narratives.

Limitations

While our PersonaMov designed for personality prediction shows superiority in most aspects, it also comes with inherent limitations.

Dialogues and character behaviors extracted from movies or TV shows may not always accurately reflect real-life personality traits due to the scripted nature of these interactions. Fictional characters are often designed to serve a narrative purpose, which might exaggerate or oversimplify certain personality traits for dramatic effect, leading to potential biases in personality prediction. The process of annotating dialogues, character relationships, and personality traits, even if partially automated, involves a degree of subjectivity. Different annotators might interpret the same dialogue or behavior differently based on their own biases and experiences, leading to inconsistencies in the dataset.

The dataset may predominantly reflect the cultural norms and values of the society in which the content was produced, potentially limiting its applicability across different cultural contexts. Our dataset is based on English movies and TV shows, so it may not interpret other non-English cultural contexts properly.

References

545

546

551

554

555

556

557

559

565

566

567

568

569

570

571

573

574

575

576

577

578

579

580

582

583

587

588

590

591

592

596

- Berne convention for the protection of literary and artistic works. https://www.wipo.int/treaties/en/ip/ berne/. Accessed [Current Date].
- Universal copyright convention. https://www.wipo. int/treaties/en/ip/ucc/. Accessed [Current Date].
- Abdullah Al Maruf, Md. Abdullah-Al Nayem, Md. Mahmudul Haque, Zakaria Masud Jiyad, Al Mamun Or Rashid, and Fahima Khanam. 2022.
 A survey on personality prediction. In *Proceedings* of the 2nd International Conference on Computing Advancements, ICCA '22, page 407–414, New York, NY, USA. Association for Computing Machinery.
- Digbalay Bose, Rajat Hebbar, Krishna Somandepalli, Haoyang Zhang, Yin Cui, Kree Cole-McLaughlin, Huisheng Wang, and Shrikanth Narayanan. 2023. Movieclip: Visual scene recognition in movies. In 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pages 2082–2091.
- Yirong Chen, Weiquan Fan, Xiaofen Xing, Jianxin Pang, Minlie Huang, Wenjing Han, Qianfeng Tie, and Xiangmin Xu. 2022. CPED: A large-scale chinese personalized and emotional dialogue dataset for conversational ai.
- Paul Costa and Robert McCrae. 2002. Personality in adulthood: A five-factor theory perspective. *Management Information Systems Quarterly MISQ*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.
- Tao Feng, Chuanyang Jin, Jingyu Liu, Kunlun Zhu, Haoqin Tu, Zirui Cheng, Guanyu Lin, and Jiaxuan You. 2024. How far are we from agi.
- Kevin A. Fischer. 2023. Reflective linguistic programming (rlp): A stepping stone in socially-aware agi (socialagi).

Hang Jiang, Xianzhe Zhang, and Jinho D. Choi. 2020. Automatic text-based personality recognition on monologues and multiparty dialogues using attentive networks and contextual embeddings (student abstract). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(10):13821–13822. 597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

- Julio C. S. Jacques Junior, Agata Lapedriza, Cristina Palmero, Xavier Baro, and Sergio Escalera. 2021. Person perception biases exposed: Revisiting the first impressions dataset. In 2021 IEEE Winter Conference on Applications of Computer Vision Workshops (WACVW). IEEE.
- Elma Kerz, Yu Qiao, Sourabh Zanwar, and Daniel Wiechmann. 2022. Pushing on personality detection from verbal behavior: A transformer meets text contours of psycholinguistic features.
- Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg. 2018. TVQA: Localized, compositional video question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Bruno Lepri, Jacopo Staiano, Giulio Rigato, Kyriaki Kalimeri, Ailbhe Finnerty, Fabio Pianesi, Nicu Sebe, and Alex Pentland. 2012. The sociometric badges corpus: A multilevel behavioral dataset for social behavior in complex organizations. In 2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Conference on Social Computing, pages 623–628.
- Bohan Li, Jiannan Guan, Longxu Dou, Yunlong Feng, Dingzirui Wang, Yang Xu, Enbo Wang, Qiguang Chen, Bichen Wang, Xiao Xu, Yimeng Zhang, Libo Qin, Yanyan Zhao, Qingfu Zhu, and Wanxiang Che. Can large language models understand you better? an MBTI personality detection dataset aligned with population traits.
- Zheng Lian, Licai Sun, Yong Ren, Hao Gu, Haiyang Sun, Lan Chen, Bin Liu, and Jianhua Tao. 2024. Merbench: A unified evaluation benchmark for multimodal emotion recognition.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach.
- OpenAI. 2023. Gpt-3.5-turbo. https://www.openai. com/research/gpt-3-5. Accessed: 2024-06-03.

OpenAI. 2024. Gpt-4-0125. GPT-4-0125 model.

- Sayantan Pal, Souvik Das, and Rohini K. Srihari. Beyond discrete personas: Personality modeling through journal intensive conversations.
- Cristina Palmero, Javier Selva, Sorina Smeureanu, Julio C. S. Jacques Junior, Albert Clapés, Alexa Moseguí, Zejian Zhang, David Gallardo, Georgina

- 655

- 660

674

682

686

695

703

Guilera, David Leiva, and Sergio Escalera. 2021. Context-aware personality inference in dyadic scenarios: Introducing the udiva dataset. In 2021 IEEE Winter Conference on Applications of Computer Vision Workshops (WACVW), pages 1-12.

- Mingwei Sun and Kunpeng Zhang. 2023. Multimodal co-attention transformer for video-based personality understanding. In 2023 IEEE International Conference on Big Data (BigData), pages 1450-1459.
- Chiu-yu Tseng, Chao-yu Su, and Tanya Visceglia. 2013. Levels of lexical stress contrast in english and their realization by 11 and 12 speak-In 2013 International Conference Oriental ers. COCOSDA held jointly with 2013 Conference on Asian Spoken Language Research and Evaluation (O-COCOSDA/CASLRE), pages 1-5.
- 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.
- Feifan Yang, Xiaojun Quan, Yunyi Yang, and Jianxing Yu. Multi-document transformer for personality detection. 35(16):14221-14229.
- Tao Yang, Jinghao Deng, Xiaojun Quan, and Qifan Wang. 2023. Orders are unwanted: Dynamic deep graph convolutional network for personality detection.
- Yangfu Zhu, Linmei Hu, Xinkai Ge, Wanrong Peng, and Bin Wu. 2022. Contrastive graph transformer network for personality detection. In Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22, pages 4559-4565. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Yaochen Zhu, Xiangqing Shen, and Rui Xia. 2023. Personality-aware human-centric multimodal reasoning: A new task.

A Definitions of Personality Models

• Myers-Briggs Type Indicator (MBTI): The MBTI categorizes personality into four dimensions. Extraversion (E) vs. Introversion (I): Extraverts are outgoing and energized by social interactions, while Introverts are reserved and energized by solitude. Sensing (S) vs. Intuition (N): Sensors focus on present, concrete information, valuing practicality, whereas Intuitives are imaginative and future-oriented, valuing abstract ideas. Thinking (T) vs. Feeling (F): Thinkers base decisions on logic and fairness, prioritizing objectivity, while Feelers base decisions on personal values and the impact on others, pri-704 oritizing harmony. Judging (J) vs. Perceiv-705 ing (P): Judgers prefer structured and orga-706 nized lives, liking plans and decisiveness, 707 while Perceivers prefer flexibility and spon-708 taneity, liking to keep their options open. 709 Each MBTI type is defined by a combina-710 tion of four cognitive functions, which can 711 be either introverted (i) or extraverted (e). 712 Extraverted Sensing (Se): Focuses on the 713 present moment and physical reality, highly 714 attuned to sensory experiences. Introverted 715 Sensing (Si): Relies on past experiences and 716 memories, valuing tradition and consistency. 717 Extraverted Intuition (Ne): Sees patterns and 718 connections, focusing on future possibili-719 ties and abstract ideas. Introverted Intuition 720 (Ni): Focuses on internal insights and fore-721 sight, seeing underlying meanings and fu-722 ture potentials. Extraverted Thinking (Te): 723 Organizes and structures the external world, 724 prioritizing logic and efficiency. Introverted 725 Thinking (Ti): Analyzes and categorizes in-726 formation internally, valuing logical consis-727 tency and understanding. Extraverted Feeling 728 (Fe): Prioritizes harmony and social values, 729 focusing on the needs and feelings of others. 730 Introverted Feeling (Fi): Values personal be-731 liefs and feelings, making decisions based on 732 inner values and ethics. 733

• Big Five Personality Traits: The Big Five model describes personality using five broad traits. Openness to Experience: High openness involves imagination and insight, while low openness involves practicality and routine. Conscientiousness: High conscientiousness is characterized by organization and dependability, while low conscientiousness is characterized by spontaneity and flexibility. Extraversion: High extraversion includes sociability and assertiveness, while low extraversion (introversion) includes reserve and solitude. Agreeableness: High agreeableness involves trust and altruism, while low agreeableness involves skepticism and competition. Neuroticism: High neuroticism involves emotional instability and anxiety, while low neuroticism involves emotional stability and calmness.

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

• Enneagram: The Enneagram classifies per-

sonality into nine types, each representing 754 different motivations and fears. Type 1: The 755 Reformer, driven by a need for perfection. Type 2: The Helper, driven by a need to be loved. Type 3: The Achiever, driven by a need for success. Type 4: The Individualist, driven by a need for uniqueness. Type 5: The In-760 vestigator, driven by a need for knowledge. Type 6: The Loyalist, driven by a need for security. Type 7: The Enthusiast, driven by a need for variety and fun. Type 8: The Challenger, driven by a need for control. Type 9: 765 The Peacemaker, driven by a need for harmony. A 2w3 individual is likely to be more 767 ambitious, charming, and goal-oriented than a typical Type 2. They still seek to help others but are also motivated by a desire for success and recognition.

• Instinctual Variants: The Instinctual Vari-772 ants theory describes three primary in-773 stinctual drives influencing behavior. Self-774 Preservation (SP): Focuses on safety, health, and comfort. Social (SO): Focuses on relationships, status, and community. Sexual 777 (SX): Focuses on intimacy, attraction, and 778 one-on-one connections. For instance, an 8w7 with a Sexual variant, is highly charismatic and seeks intense and passionate connections with others. He or she is bold and assertive, often focusing his or her energy on 783 building strong, impactful relationships.

B Data Alignment Algorithm

786

788

792

794

797

798

799

802

The details of data alignment algorithm are as follows:

 Preprocess the raw data Firstly, we divide the scripts into several scenes according to the coherence in language of camera, instead of randomly clipping in a certain time period. This segmentation is guided by explicit scene transition cues found in movie scripts, such as "CUT TO:" or scene location indicators. For TV show scripts, which might lack uniform scene transition markers, we identify scene changes by detecting pauses exceeding 3 seconds between utterances.

2. *Match the utterance* This algorithm is rooted in the comparison of utterances from original scripts and subtitles based on a similarity threshold. If the similarity between a pair

Algorithm 1: Scripts and Subtitles Matching

Input: Script, Subtitles
Output: Updated subtitles with speaker names
1: $dial\&speakers \leftarrow empty$
2: threshold $\leftarrow 0.8$
3: for scene in Script do
4: for <i>Dials</i> in scene do
5: Extract <i>speaker</i> and <i>dial</i> from <i>Dials</i>
6: $dial\&speakers \leftarrow speaker, dial$
7: end for
8: end for
9: for <i>subtitle</i> in <i>Subtitles</i> do
10: $match_score \leftarrow 0$
11: $match_speaker \leftarrow Null$
12: for <i>line</i> in <i>subtitle</i> do
13: for speaker, dial in dial&speakers do
14: $score \leftarrow Similar(subtitle, dial)$
15: if score <i>i</i> match_score then
16: Update <i>match_score</i> and <i>match_speaker</i>
17: end if
18: end for
19: if $match_score \ge threshold$ then
20: Update <i>line</i> with <i>match_speaker</i>
21: end if
22: end for
23: Update <i>subtitle</i>
24: end for
25: return Updated Subtitles

of utterances meets or exceeds this threshold, the character's name is accurately associated with the utterance.

3. *Rematch with the slide window* Basically, the content in scripts is slightly different with the subtitles, because the director may have improvised on the set. Thus, we introduce a slide window algorithm to evaluate the utterance-level similarity. As shown in Algorithm 2, we set a window to slide over the script and, for each utterance, compare the content inside the window with each subtitle entry to get the similarity of the paragraph in the window.

815

816

803

804

Algorithm 2: Slide Window Matching

Input: Script, Subtitles Output: Updated subtitles 1: $window_size \leftarrow 10$ 2: $threshold \leftarrow 0.8$ 3: $matches \leftarrow empty_list$ 4: for $i \leftarrow 0$ to $Len(Script) - window_size$ do window5: \leftarrow slice(scriptTokens, i, i + $window_size)$ $match_score \leftarrow 0$ 6: for $j \leftarrow 0$ to Len(Subtitles) - 1 do 7: 8: $score \leftarrow Similar(window, Subtitles[j])$ 9: if score ; match_score then 10: Update match_score 11: end if 12: end for 13: if $match_score \geq threshold$ then 14: $matches \leftarrow Subtitles[j]$ 15: end if 16: end for

17: return Updated Subtitles with matches