

PersonaMovs: A Multimedia Conversational Dataset for Dynamic Personality Analysis

Anonymous ACL submission

Abstract

Automatic personality detection has evolved from simple text classification to sophisticated multimodal analysis, recognizing the multi-dimensional 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 (PM), 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. PM not only addresses the challenges of existing datasets by offering majority-voted personality annotations and detailed relations networks but also paves the way for advanced analysis of personality dynamics across various contexts.

1 Introduction

Personality is a comprehensive yet complex trait that encapsulates individual differences in patterns of thinking, feeling, and behaving (Costa and McCrae, 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). In recent years, there has been a burgeoning interest in automatic personality detection, marking a significant shift from traditional methods to innovative computational approaches. At the very beginning, owing to the limitations of multimedia model and computational power, researchers only treat personality prediction as a straightforward text classification task, aiming to decipher personality traits from the digital footprints indi-

viduals leave online (Kerz et al., 2022; Yang et al.). However, as shown in Figure 1, researchers have increasingly recognized that personality is manifested via multi-dimensions, with nuances that pure text-based analysis cannot fully capture (Al Maruf et al., 2022; Zhu et al., 2022; Bose et al., 2023). This revelation has propelled the move towards multimodal personality detection as the mainstream methodology.

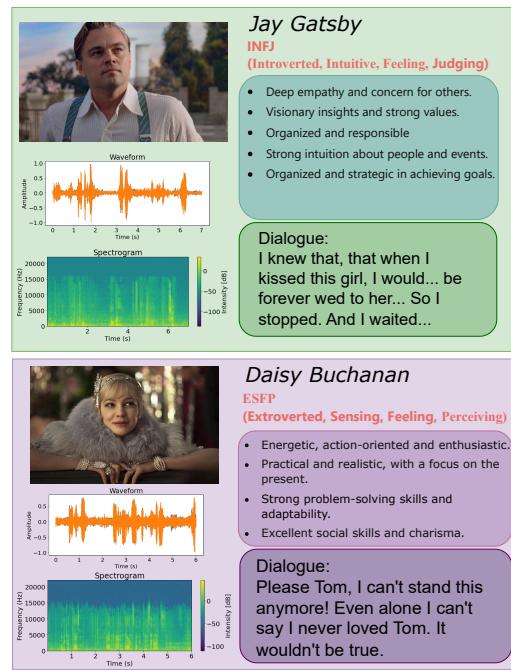


Figure 1: The Distinctive Features in Three Modalities for Personality Prediction

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 essential for developing models that accurately reflect the complexity of human personality. Naturally, a lot of multimodal datasets were released in recent years. There have been a few attempts in multimodal personality dataset construction (Palmero et al., 2021; Junior et al., 2021; Jiang et al., 2020; Chen et al., 2022). There are also many multimodal datasets used to perform other tasks, and some personality prediction works will modify their datasets to adapt the personality context. For instance, TVQA (Lei et al., 2018) is a large dataset which is initially designed to do the visual question answering task. It is used frequently in our research field because of its large scale.

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) **Dynamic Nature of Personality**: From a psychological standpoint, it is essential to recognize that personality is not a static attribute but one that evolves in response to environmental contexts (Palmero et al., 2021), which current datasets do not adequately capture.

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. As for how to get the personality dynamics, incorporating relationship networks into personality prediction models offers a solution to this issue. Such networks provide a

rich context for observing and understanding individual behaviors, preferences, and traits, reflecting the interconnectedness of personality with social and environmental factors.

Against these backdrops, we introduce the PersonaMovs (PM), a comprehensive dataset that starkly contrasts with existing offerings in several key aspects. PM 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), 7 kinds of Social Relations and 8 attitudes for the Emotion Relations. Our analysis highlights substantial quantity and diversity in content, adequate experiments on different models with all modalities and personality dynamics discovery.

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 crowdsourced personality, emotion, 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.
- For the first time, we categorize 15 types of relations to depict the dynamics of character interactions on a scene-by-scene basis, enabling a granular analysis of personality **dynamics** through relations networks. Guided by the relations networks, we identify psychological phenomena in both short-term and long-term conversational contexts, which largely explain the personality dynamics statistically.

2 Dataset Design

This paper introduces a new multimedia personality dataset, PersonaMovs, which is the largest of exist-

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

²<https://anonymous.4open.science/r/sample-of-MMPD-F26F/>

| 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

ing multimedia datasets. In this section, we provide a specific description about our dataset in terms of design principles and the structure in details.

2.1 Design Principles

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.

2.1.1 Multiple Personality Models are Needed

In constructing such a dataset for personality prediction, incorporating four distinct personality models, provides a comprehensive framework for understanding the multifaceted nature of human Personality. To this end, evolving four distinct personality models—Myers-Briggs Type Indicator (MBTI), Big Five, Enneagram, and Instinctual Variant—into our dataset construction is essential. Each of these models provides a unique lens through which to view and interpret personality traits, offering complementary insights that are critical for a holistic understanding. By integrating these four models, we aim to construct a dataset that not only captures the complexity of human personality but also facilitates nuanced predictions. This comprehensive framework acknowledges the diversity of human experience and the need for multidimensional analysis to truly understand and predict personality dynamics. The complete definitions can be found in

Appendix A.

2.1.2 Personality as Fluidity Rather than Stability

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 two main categories of relations among characters (more details in Appendix B):

1) **Social Relations:** 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. This approach acknowledges that personality is not solely a matter of internal traits and instincts, but is also fundamentally shaped and expressed through interactions with others in various domains of life.

2) **Emotion Relations:** The social relations above are relatively unchangeable, not depicting the attitudes towards someone else. So we define another 8 types for the emotion relations, as the aid for the comprehension of personality. We choose *fondness*, *jealousy*, *aversion*, *pity*, *respect*, *hostility*, *envy* and *gratitude* as our annotators for the emotional relations, which concludes the diverse attitudes in human’s daily life. To annotate these two relations, we select a binary tuple (e.g., *Gil and Adriana: (Romantic, Fondness)*) to represent the relations combination for each pair of characters based on different scenes.

2.2 Structure of PersonaMovs

Aiming to deliver a tidy and readable structure, there is no more suitable file types than JSON format. We distribute different scenes in a single JSON file with index. For each movie or TV show, the video clips with the corresponding JSON and audio files are stored in the same directory and each data point in this dataset is centered around a dialogue scene that involves several interlocutors

235 based on the original scripts.
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251

As shown in Figure 2, each video clip of our dataset is tagged with a “Scene” identifier, which likely refers to a specific segment or moment within a larger narrative or dataset. The “Dialogue” field contains an array of objects, each providing a detailed description of a scene and dialogues between characters. The dialogues are presented with time corresponding timestamps, too. The “Relationship” field within this object provides a summary or interpretation of their interaction, in this case, indicating a professional relationship with an element of fondness between Gatsby and Daisy. Finally, the “Personality” section provides personality profiles for the characters mentioned in the scene, where their personality type distribution are listed along with a “Votes Distribution” field.

252 3 Methodology

253 In this section, we outline the methodology em-
254 ployed to gather, process, and annotate the data for
255 our study. We begin by detailing the sources of
256 our multimodal video data and personality labels,
257 focusing on how we efficiently align subtitles with
258 original scripts to ensure accurate temporal and
259 character associations. And we also present our
260 annotation process, explaining how we leverage
261 the ChatGPT API to automatically annotate social
262 and emotional relations among characters within
263 the text data.

264 3.1 Source of Data

265 Our data source contains mainly two parts, the
266 multimodal video data and personality labels. For
267 video data, we include 14 different genres of TV
268 series and movies via an open-source website¹, and
269 for the scripts and subtitles, we also find other open-
270 source websites²³ for research offering the free
271 scripts and subtitles of many famous movie and
272 television programs. Considering the insufficient
273 labeling method of existing works, we collect the
274 personality annotations from personality database
275 website as well as the voting distribution and align
276 them to correctly scripts.

277 3.2 Data Alignment Process

278 As subtitle contain temporal information and origi-
279 nal scripts associate utterances with characters, we

280 are supposed to align them properly as efficient
281 as possible. However, most of the existing multi-
282 modal datasets annotate the timestamps manually
283 with taking up a great deal of time. There are
284 also some works which utilize different automatic
285 tools to align the utterances with their correspond-
286 ing information. For instance, Lian et al. (2024)
287 use an Automatic Sound Recognition (ASR) tool
288 called Gentle⁴ to get the timestamps for the utter-
289 ances. To streamline the process of aligning dia-
290 logue utterances with their respective timestamps
291 and speakers from subtitles, we propose an efficient
292 method leveraging a fuzzy matching algorithm (see
293 Appendix C). Following successful alignment, we
294 proceed to segment the video content into distinct
295 scenes according to the timestamps. Besides, we
296 use FFmpeg¹ to extract the audio track from the
297 video clips and output it as a .mp3 file.

298 3.3 Annotation Process

299 We construct a process to automatically annotate
300 the social and emotion relations among characters
301 by using ChatGPT API (OpenAI, 2023). Only text
302 data are supposed to be processed, thus we choose
303 *gpt-3.5-turbo-1106* pre-trained model to annotate
304 our dataset. Since we preprocess the text data and
305 divide them into scenes, we design a prompt to ask
306 ChatGPT for identifying both social and emotion
307 relations for every single scene. When annotating
308 our data, we encounter a challenge in representing
309 unidirectional affectionate relationships, where A
310 likes B, but B does not reciprocate the feelings.
311 While social relations do not present this issue,
312 emotion relations require a solution to capture this
313 directional information. We address this by inter-
314 pretting the relative position of different characters
315 in the tuple. For example, A and B (family, fond-
316 ness) indicates that A has positive feelings towards
317 B. Conversely, if B has positive feelings towards
318 A, the tuple would be B and A, (family, fondness).
319 This approach allows us to clearly represent the
320 directionality of emotional relations.

321 Based on the definitions of relations, we design
322 this prompt for relations annotation (Fig 4). The
323 prompt categorizes relations into seven social and
324 eight emotion types, ensuring comprehensive cov-
325 erage of human interactions. Note that we require
326 ChatGPT to generate the responses following our
327 format strictly so that we could better manipulate
328 them flexibly.

¹<https://yts.mx/>

²<https://www.simplyscripts.com/>

³<https://subscene.com/>

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

¹<https://ffmpeg.org/>

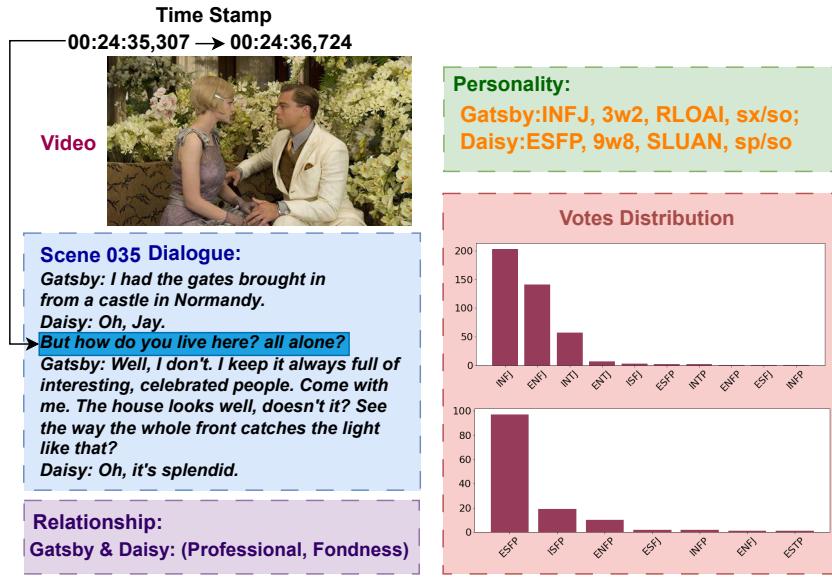


Figure 2: A sample from PersonaMovs

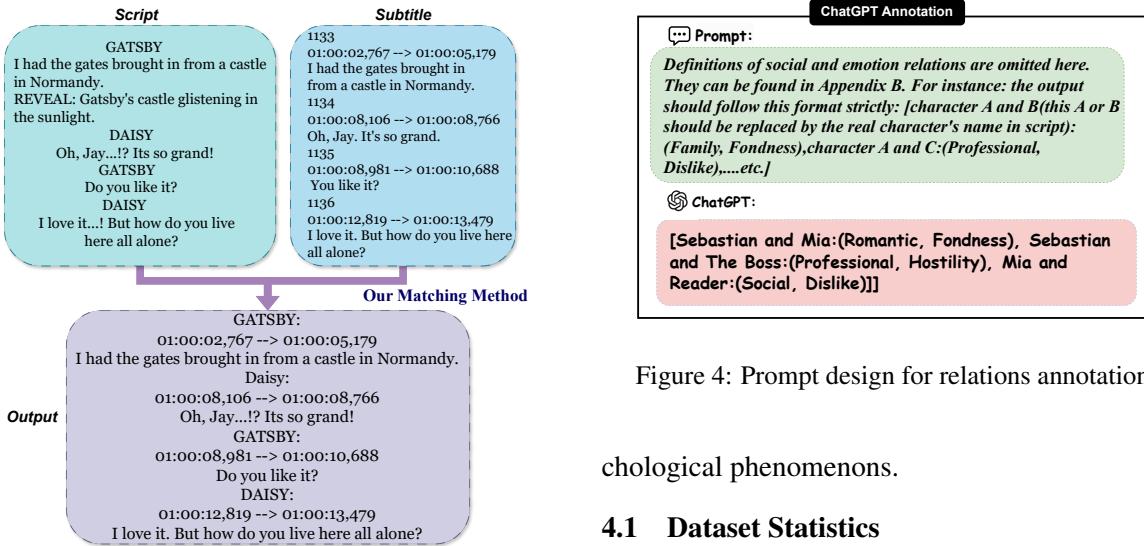


Figure 3: Process of data alignment

4 Evaluation

We present the basic statistics of our dataset in the first part, and then we evaluate the accuracy 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 relations annotation. To discover interesting topics about personality dynamics, we focus on those changes of personality and discover several interesting psy-

chological phenomena.

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 two types of relations, which indicates the diversity in terms of interaction scenarios.

4.1.1 Algorithm Evaluation

To evaluate the performance of our character-to-subtitle matching algorithm, we randomly sample a test case comprising over 50 dialogues and 600 utterances from a variety of genres, including 10 films and TV series. We manually check

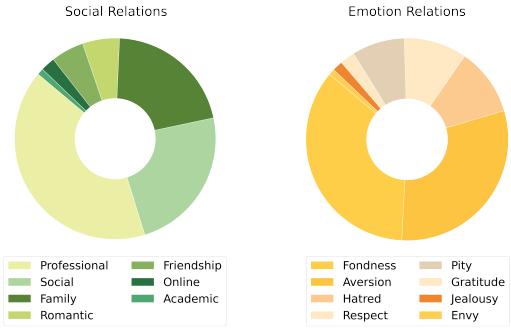


Figure 5: Distribution of social and emotion relations

the aligned characters’ name based on the script. Our primary metrics for assessment is 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 2: Accuracy and running time per dialogue of subtitle matching algorithm

4.1.2 Annotation Accuracy

Using ChatGPT to annotate relations for the characters is not a completely worthwhile method. To measure the automatic annotation accuracy, we sampled 235 scenes randomly and involved 5 human labelers on relations annotation. These labelers are in their mid-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 and emotion relations after aligned video. We continue to compare the automatically annotated results to the human-labeled ground truth. The outcome shows that both social and emotional relationship annotations are dependable, with the accuracy reaching 95% and 84% respectively.

| Task | Movies | TV | Total |
|-------------------|--------|--------|--------|
| Social Relations | 98.21% | 93.91% | 95.78% |
| Emotion Relations | 82.04% | 84.46% | 84.01% |

Table 3: Accuracy of relations annotation.

The dataset’s foundation on crowdsourced voting allows for an in-depth analysis of subjective 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.

4.2 Experiment Results

4.2.1 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), AttRCNN (Xue et al., 2018), GPT-3.5 (OpenAI, 2023), GPT-4 (OpenAI, 2024) and MCT (Sun and Zhang, 2023). Table 4 shows the accuracy of our dataset is apparently lower than other competing datasets. One of the main challenges we observed was the complexity and diversity of our dataset compared to other multimedia 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 4: 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 terms of personality detection: 1) Our dataset captures 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 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 successful. This helps in distinguishing truly advanced models from those that perform well only on simpler

421 tasks. Additionally, one notable observation from
 422 the results is that the MCT model, which leverages
 423 three modalities (text, audio, and video), does not
 424 outperform the GPT-4 model, which uses only two
 425 modalities (text and video). This performance gap
 426 suggests that Large Language Model outperforms
 427 the small model on this task, even though the latter
 428 uses more modalities.

429 4.2.2 The Importance of Multi-Modality

430 We conducted a series of ablation experiments to as-
 431 sess the impact of different modalities and relations
 432 annotations on the performance of personality pre-
 433 diction models. The experiments were designed to
 434 understand how the exclusion of specific modalities
 435 or relations annotations affects the overall predic-
 436 tion accuracy.

| Method | Modality | Accuracy |
|------------|----------|----------|
| MCT | T & A | 66.13 |
| | T & V | 67.91 |
| | T only | 63.43 |
| | T, A & V | 68.47 |
| GPT-4-0125 | T only | 70.20 |
| | T & V | 76.90 |

437 Table 5: Ablation experiment on different modalities

438 Table 5 presents the results of ablation experi-
 439 ments where different combinations of video and
 440 audio modalities were excluded. The result un-
 441 derscore the critical importance of using multiple
 442 modalities to achieve higher accuracy in per-
 443 sonality prediction tasks. Models that leverage both
 444 audio and video data, in addition to text, consist-
 445 ently outperform those that rely solely on textual
 446 data.

| Method | With Relations | Without Relations |
|------------|----------------|-------------------|
| BERT | 53.88 | 52.94 |
| Roberta | 59.21 | 58.39 |
| GPT-4-0125 | 73.22 | 70.20 |

447 Table 6: Ablation experiment on relations annotations.

448 Table 6 shows the results of ablation experiments
 449 focusing on the inclusion or exclusion of relations'
 450 annotations which finds the relations annotations
 451 tend to slightly enhance the performance. This
 452 highlights the importance of including rich context-
 453 ual information to improve the accuracy of person-
 454 ality prediction models.

455 The multimodal nature of the dataset (incorporat-
 456 ing video, audio, textual, and crowd-sourced data)
 457 enables comprehensive studies that integrate dif-
 458 ferent data types to understand personality. This

459 could lead to the development of new theories or
 460 the refinement of existing ones that account for the
 461 complexity of personality as depicted through vari-
 462 ous media. It could also foster interdisciplinary re-
 463 search, combining insights from psychology, com-
 464 puter science, linguistics, and media studies.

465 4.2.3 Personality Dynamics

466 Movies and TV series and their characters often
 467 evolve over time, offering a fertile ground for study-
 468 ing personality dynamics. The dataset allows for
 469 longitudinal studies on how characters' personali-
 470 ties change in response to narrative events, rela-
 471 tionships, and challenges. This could lead to new
 472 models that explain personality development and
 473 dynamics in complex social settings, bridging nar-
 474 rative theory and psychological research.

475 We manually select two famous characters to
 476 find their potential change of personality, and we
 477 discover there are two types of personality shifts.
 478 In the short term, people show different personali-
 479 ties depending on their relations with interlocutor.
 480 For example, as shown in Figure 6, Jay Gatsby
 481 from the famous romantic movie called “*The Great*
 482 *Gatsby*” behaves as an INFJ in front of his beloved
 483 Daisy and as an ENFJ in front of his business part-
 484 ners. In the long term, people may change their
 485 personality due to major turning point of life. Like
 486 Mia Dolan in “*La La Land*”, she was always ENFJ,
 487 but after the breakup she became an ISFJ. The pre-
 488 diction results generated by GPT-4 align with the
 489 peaks in the voting distribution, indicating that this
 490 personality shift is observable within our realistic
 491 dataset.

492 According to our finding, we conduct statisti-
 493 cal analysis based on our dataset to figure out if
 494 there exists certain personalities that are easily at-
 495 tracted to each other. To analyze the patterns of
 496 personality attraction, we focus on identifying pairs
 497 of personalities that frequently appear together bi-
 498 directionally in fondness, aversion, romantic and
 499 friendship relations. Figure 7 presents the favorite
 500 network with 16 MBTI personality types, providing
 501 a clear visual summary of statistical findings. The
 502 size of each node is proportional to the number of
 503 connections (degree) it has, which means personal-
 504 ity types with more relationships are represented by
 505 larger nodes. The color of the edges represents the
 506 weight of the relationship between the personality
 507 types. Darker edges indicate a higher frequency or
 508 stronger relationship. Based on these vivid figures,
 509 we can discover very interesting psychological phe-

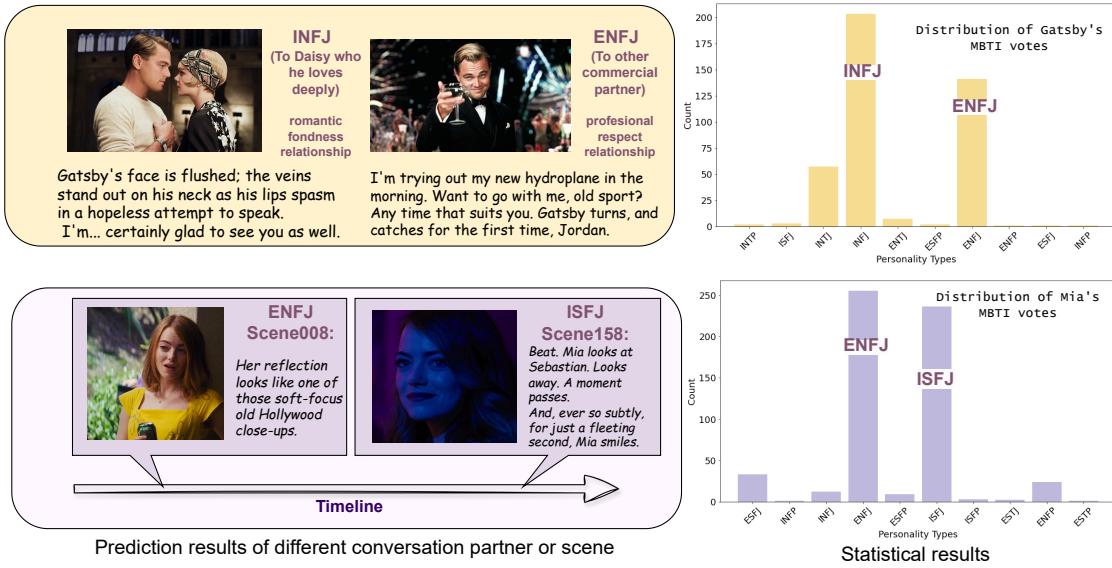


Figure 6: Case study for personality dynamics.

nomenons. For instance, ISTP is the most popular personality since almost every other personality has a fondness relation with it, and ESTP prefers to be around ESFP, ENFP and people with the same personality as themselves. ESFP may not like people with same personality because it has a dark circle on itself.

ability through interactions. This aspect can support research into how different personality types influence and are influenced by social networks, both within narrative contexts and long-term conversions. It provides a basis for computational models that simulate personality dynamics in social networks, potentially informing theories on social behavior, conflict resolution, and group dynamics.

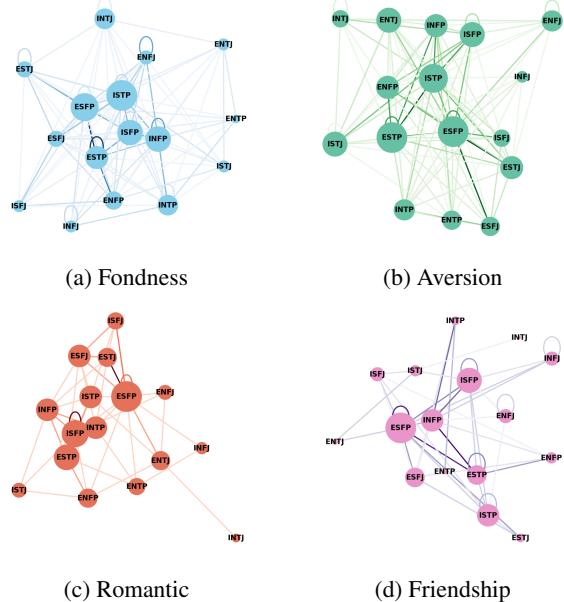


Figure 7: Favorite Networks with Different Personalities about Four Relations

By including data on social and emotion relations between characters, the dataset opens new pathways for exploring the dynamics of person-

5 Conclusion

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

Copyright Concerns

Copyright © [2024] by the authors. The movies and TV series included in this dataset are copyrighted by their respective copyright owners and are used in this work for academic and research purposes under fair use guidelines or specific per-

547 missions obtained from the copyright holders. This
548 does not imply endorsement by or affiliation with
549 the copyright owners. Use of these materials is
550 limited to the scope of the permission granted and
551 is not intended for commercial distribution.

552 Limitations

553 While our PersonaMosaic designed for personality
554 prediction shows superiority in most aspects, it also
555 comes with inherent limitations.

556 Dialogues and character behaviors extracted
557 from movies or TV shows may not always accu-
558 rately reflect real-life personality traits due to the
559 scripted nature of these interactions. Fictional char-
560 acters are often designed to serve a narrative pur-
561 pose, which might exaggerate or oversimplify cer-
562 tain personality traits for dramatic effect, leading
563 to potential biases in personality prediction.

564 The process of annotating dialogues, character
565 relationships, and personality traits, even if par-
566 tially automated, involves a degree of subjectivity.
567 Different annotators might interpret the same dia-
568 logue or behavior differently based on their own
569 biases and experiences, leading to inconsistencies
570 in the dataset.

571 The dataset may predominantly reflect the cul-
572 tural norms and values of the society in which the
573 content was produced, potentially limiting its ap-
574 plicability across different cultural contexts. Our
575 dataset is based on English movies and TV shows,
576 so it may not interpret other non-English cultural
577 contexts properly.

578 References

579 Abdullah Al Maruf, Md. Abdullah-Al Nayem,
580 Md. Mahmudul Haque, Zakaria Masud Jiyad,
581 Al Mamun Or Rashid, and Fahima Khanam. 2022.
582 [A survey on personality prediction](#). In *Proceedings
583 of the 2nd International Conference on Computing
584 Advancements, ICCA '22*, page 407–414, New York,
585 NY, USA. Association for Computing Machinery.

586 Digbalay Bose, Rajat Hebbar, Krishna Somandepalli,
587 Haoyang Zhang, Yin Cui, Kree Cole-McLaughlin,
588 Huisheng Wang, and Shrikanth Narayanan. 2023.
589 [Movieclip: Visual scene recognition in movies](#). In
590 *2023 IEEE/CVF Winter Conference on Applications
591 of Computer Vision (WACV)*, pages 2082–2091.

592 Yirong Chen, Weiquan Fan, Xiaofen Xing, Jianxin
593 Pang, Minlie Huang, Wenjing Han, Qianfeng Tie,
594 and Xiangmin Xu. 2022. [CPED: A large-scale chi-](#)
595 [nese personalized and emotional dialogue dataset for
596 conversational ai](#).

597 W. Andrew Collins and L. Alan Sroufe. 1999. *Ca-*
598 *capacity for Intimate Relationships: A Developmental*
599 *Construction*, Cambridge Studies in Social and Emo-
600 *tional Development*, page 125–147. Cambridge Uni-
601 *versity Press.*

602 Paul Costa and Robert McCrae. 2002. [Personality in
603 adulthood: A five-factor theory perspective](#). *Man-*
604 *agement Information Systems Quarterly - MISQ*.

605 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
606 Kristina Toutanova. 2019. [Bert: Pre-training of deep
607 bidirectional transformers for language understand-](#)
608 [ing](#).

609 Robert A. Emmons and Michael E. McCullough. 2004.
610 [The Psychology of Gratitude](#). Oxford University
611 Press.

612 Tao Feng, Chuanyang Jin, Jingyu Liu, Kunlun Zhu,
613 Haoqin Tu, Zirui Cheng, Guanyu Lin, and Jiaxuan
614 You. 2024. [How far are we from agi](#).

615 Kevin A. Fischer. 2023. [Reflective linguistic pro-](#)
616 [gramming \(rlp\): A stepping stone in socially-aware agi](#)
617 [\(socialagi\)](#).

618 Hang Jiang, Xianzhe Zhang, and Jinho D. Choi. 2020.
619 [Automatic text-based personality recognition on](#)
620 [monologues and multiparty dialogues using atten-](#)
621 [tive networks and contextual embeddings \(student](#)
622 [abstract\)](#). *Proceedings of the AAAI Conference on*
623 *Artificial Intelligence*, 34(10):13821–13822.

624 Julio C. S. Jacques Junior, Agata Lapedriza, Cristina
625 Palmero, Xavier Baro, and Sergio Escalera. 2021.
626 [Person perception biases exposed: Revisiting the](#)
627 [first impressions dataset](#). In *2021 IEEE Winter Con-*
628 *ference on Applications of Computer Vision Work-*
629 *shops (WACVW)*. IEEE.

630 Elma Kerz, Yu Qiao, Sourabh Zanwar, and Daniel
631 Wiechmann. 2022. [Pushing on personality detection](#)
632 [from verbal behavior: A transformer meets text con-](#)
633 [tours of psycholinguistic features](#).

634 Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg.
635 2018. [TVQA: Localized, compositional video ques-](#)
636 [tion answering](#). In *Proceedings of the 2018 Con-*
637 *ference on Empirical Methods in Natural Language*
638 *Processing*. Association for Computational Linguis-
639 *tics*.

640 Bruno Lepri, Jacopo Staiano, Giulio Rigato, Kyriaki
641 Kalimeri, Ailbhe Finnerty, Fabio Pianesi, Nicu Sebe,
642 and Alex Pentland. 2012. [The sociometric badges](#)
643 [corpus: A multilevel behavioral dataset for social](#)
644 [behavior in complex organizations](#). In *2012 Inter-*
645 *national Conference on Privacy, Security, Risk and*
646 *Trust and 2012 International Conference on Social*
647 *Computing*, pages 623–628.

648 Zheng Lian, Licai Sun, Yong Ren, Hao Gu, Haiyang
649 Sun, Lan Chen, Bin Liu, and Jianhua Tao. 2024.
650 [Merbench: A unified evaluation benchmark for mul-](#)
651 [timodal emotion recognition](#).

- 652 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
 653 dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
 654 Luke Zettlemoyer, and Veselin Stoyanov. 2019.
 655 **Roberta: A robustly optimized bert pretraining ap-**
 656 **proach.**
- 657 OpenAI. 2023. Gpt-3.5-turbo. <https://www.openai.com/research/gpt-3-5>. Accessed: 2024-06-03.
- 658
- 659 OpenAI. 2024. **Gpt-4-0125**. GPT-4-0125 model.
- 660 Cristina Palmero, Javier Selva, Sorina Smeureanu,
 661 Julio C. S. Jacques Junior, Albert Clapés, Alexa
 662 Moseguí, Zejian Zhang, David Gallardo, Georgina
 663 Guilera, David Leiva, and Sergio Escalera. 2021.
 664 **Context-aware personality inference in dyadic sce-**
 665 **narios: Introducing the udva dataset.** In *2021 IEEE*
 666 *Winter Conference on Applications of Computer Vi-*
 667 *sion Workshops (WACVW)*, pages 1–12.
- 668 Mingwei Sun and Kunpeng Zhang. 2023. **Multimodal**
 669 **co-attention transformer for video-based personality**
 670 **understanding.** In *2023 IEEE International Confer-*
 671 *ence on Big Data (BigData)*, pages 1450–1459.
- 672 Chiu-yu Tseng, Chao-yu Su, and Tanya Visceglia.
 673 2013. **Levels of lexical stress contrast in en-**
 674 **glish and their realization by 11 and 12 speak-**
 675 **ers.** In *2013 International Conference Oriental*
 676 *COCOSDA held jointly with 2013 Conference on*
 677 *Asian Spoken Language Research and Evaluation*
 678 (*O-COCOSDA/CASLRE*), pages 1–5.
- 679 Di Xue, Lifa Wu, Zheng Hong, Shize Guo, Liang Gao,
 680 Zhiyong Wu, Xiaofeng Zhong, and Jianshan Sun.
 681 2018. **Deep learning-based personality recognition**
 682 **from text posts of online social networks.** *Applied*
 683 *Intelligence*, 48.
- 684 Feifan Yang, Xiaojun Quan, Yunyi Yang, and Jianxing
 685 **Yu. Multi-document transformer for personality de-**
 686 **tection.** 35(16):14221–14229.
- 687 Tao Yang, Jinghao Deng, Xiaojun Quan, and Qifan
 688 **Wang. 2023. Orders are unwanted: Dynamic deep**
 689 **graph convolutional network for personality detec-**
 690 **tion.**
- 691 Yangfu Zhu, Linmei Hu, Xinkai Ge, Wanrong Peng,
 692 and Bin Wu. 2022. **Contrastive graph transformer**
 693 **network for personality detection.** In *Proceedings*
 694 *of the Thirty-First International Joint Conference on*
 695 *Artificial Intelligence, IJCAI-22*, pages 4559–4565.
 696 International Joint Conferences on Artificial Intelli-
 697 *gence Organization. Main Track.*
- 698 Yaochen Zhu, Xiangqing Shen, and Rui Xia. 2023.
 699 **Personality-aware human-centric multimodal rea-**
 700 **soning: A new task.**

701 A Definitions of Personality Models

- 702 • **Myers–Briggs Type Indicator (MBTI):** The
 703 MBTI categorizes personality into four dimen-
 704 sions. Extraversion (E) vs. Introversion (I):

705 Extraverts are outgoing and energized by so-
 706 cial interactions, while Introverts are reserved
 707 and energized by solitude. Sensing (S) vs. In-
 708 tuitives focus on present, concrete
 709 information, valuing practicality, whereas In-
 710 tuitives are imaginative and future-oriented,
 711 valuing abstract ideas. Thinking (T) vs. Feel-
 712 ing (F): Thinkers base decisions on logic and
 713 fairness, prioritizing objectivity, while Feelers
 714 base decisions on personal values and the im-
 715 pact on others, prioritizing harmony. Judging
 716 (J) vs. Perceiving (P): Judgers prefer struc-
 717 tured and organized lives, liking plans and de-
 718 cisiveness, while Perceivers prefer flexibility
 719 and spontaneity, liking to keep their options
 720 open. Each MBTI type is defined by a combi-
 721 nation of four cognitive functions, which can
 722 be either introverted (i) or extraverted (e). Ex-
 723 traverted Sensing (Se): Focuses on the present
 724 moment and physical reality, highly attuned
 725 to sensory experiences. Introverted Sensing
 726 (Si): Relies on past experiences and memories,
 727 valuing tradition and consistency. Extraverted
 728 Intuition (Ne): Sees patterns and connections,
 729 focusing on future possibilities and abstract
 730 ideas. Introverted Intuition (Ni): Focuses on
 731 internal insights and foresight, seeing under-
 732 lying meanings and future potentials. Ex-
 733 traverted Thinking (Te): Organizes and struc-
 734 tures the external world, prioritizing logic and
 735 efficiency. Introverted Thinking (Ti): Ana-
 736 lyzes and categorizes information internally,
 737 valuing logical consistency and understanding.
 738 Extraverted Feeling (Fe): Prioritizes harmony
 739 and social values, focusing on the needs and
 740 feelings of others. Introverted Feeling (Fi):
 741 Values personal beliefs and feelings, making
 742 decisions based on inner values and ethics.

- **Big Five Personality Traits:** The Big Five
 743 model describes personality using five broad
 744 traits. Openness to Experience: High open-
 745 ness involves imagination and insight, while
 746 low openness involves practicality and rou-
 747 tine. Conscientiousness: High conscientious-
 748 ness is characterized by organization and de-
 749 pendability, while low conscientiousness is
 750 characterized by spontaneity and flexibility.
 751 Extraversion: High extraversion includes so-
 752 ciability and assertiveness, while low extraversion
 753 (introversion) includes reserve and soli-
 754 tude. Agreeableness: High agreeableness in-

| 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 7: Descriptions of Social Relations

| Relations type | Description |
|----------------|--|
| Fondness | A positive emotion characterized by a person's fondness for another. |
| Jealousy | Unhappy and angry because someone has something that you want. |
| Aversion | A negative emotion, referring to a feeling of disfavor towards someone. |
| Pity | A feeling of sadness for someone else's difficult situation. |
| Respect | Admiration felt or shown for someone that you believe has good ideas or qualities. |
| Hostility | An unfriendly or unkindness towards someone or something. |
| Envy | A discontented feeling when a person desires what someone else has. |
| Gratitude | An emotion of being thankful for someone else's help or kind actions. |

Table 8: Description of Emotion Relations

volves 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.

- **Enneagram:** The Enneagram classifies personality into nine types, each representing different motivations and fears. Type 1: The 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 Investigator, 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: The Peacemaker, driven by a need for harmony. A 2w3 individual is likely to be more 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 Variants theory describes three primary instinctual drives influencing behavior. Self-Preservation (SP): Focuses on safety, health, and comfort. Social (SO): Focuses on relationships, status, and community. Sexual (SX): Focuses on intimacy, attraction, and 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 building strong, impactful relationships.

B Definitions of Relations

Human social networks are complex and multi-faceted. By categorizing relations, we can better understand the dynamics and 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. Table 7 and Table 8 provide a structured approach to understanding the complex web of relations that individuals navigate. By categorizing these relations into social and emotional types, we can better analyze and predict personality dynamics in various contexts (Collins and Sroufe, 1999; Emmons and McCullough, 2004).

C Data Alignment Algorithm

The details of data alignment algorithm are as follows:

1. *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

756
757
758
759
760
761

762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779

780
781
782
783
784
785
786
787

788
789
790
791
792

793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808

809
810
811
812
813
814
815
816
817

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 Dials in scene do
5:     Extract speaker and dial from Dials
6:     dial&speakers  $\leftarrow$  speaker, dial
7:   end for
8: end for
9: for subtitle in Subtitles do
10:   match_score  $\leftarrow$  0
11:   match_speaker  $\leftarrow$  Null
12:   for line in subtitle do
13:     for speaker, dial in dial&speakers do
14:       score  $\leftarrow$  Similar(subtitle, dial)
15:       if score > match_score then
16:         Update match_score and match_speaker
17:       end if
18:     end for
19:     if match_score  $\geq$  threshold then
20:       Update line with match_speaker
21:     end if
22:   end for
23:   Update subtitle
24: end for
25: return Updated Subtitles
```

818 “CUT TO:” or scene location indicators. For
819 TV show scripts, which might lack uniform
820 scene transition markers, we identify scene
821 changes by detecting pauses exceeding 3 sec-
822 onds between utterances.

- 823 2. *Match the utterance* This algorithm is rooted
824 in the comparison of utterances from origi-
825 nal scripts and subtitles based on a similarity
826 threshold. If the similarity between a pair of
827 utterances meets or exceeds this threshold, the
828 character’s name is accurately associated with
829 the utterance.
- 830 3. *Rematch with the slide window* Basically, the
831 content in scripts is slightly different with the
832 subtitles, because the director may have im-
833 provised on the set. Thus, we introduce a slide
834 window algorithm to evaluate the utterance-
835 level similarity. As shown in Algorithm 2,
836 we set a window to slide over the script and,
837 for each utterance, compare the content inside
838 the window with each subtitle entry to get the
839 similarity of the paragraph in the window.

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
5:   window  $\leftarrow$  slice(scriptTokens, i, i + window_size)
6:   match_score  $\leftarrow$  0
7:   for j  $\leftarrow$  0 to Len(Subtitles) - 1 do
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
```
