

Investigating the Effects of Big Five Personality Traits on Suicide Severity Level Detection with Large Language Models

Anonymous ACL submission

Abstract

Sensitive content warning: This paper contains sensitive content related to suicide.

Suicide is one of the leading causes of global mortality, making risk detection a critical public health priority. Although psychological studies have observed a tangible link between personality and suicide, this relationship has yet to be empirically tested and applied using LLMs. In this study, we investigate the ability of LLMs to observe and leverage the relation between these two domains during suicide severity level detection. We propose inducing user persona via profiles compiled using the Big Five personality traits within a zero-shot setting, assessing the impact of persona information on detection performance. Experimental results demonstrated that the persona induction approach showed a marginal but positive impact on detection for most models, suggesting LLMs partially comprehend the inter-domain relationship and that further refinement could significantly boost performance.

1 Introduction

Suicide remains a significant cause of death worldwide, and detection of suicidal risk for preventive intervention remains a critical challenge in global public health (Ghosh et al., 2020).

A tangible link between personality and the manifestation of suicidal ideations and behaviours has been consistently demonstrated by findings in psychology. Personality, which comprises the unique combination of characteristics that form an individual’s distinctive character. The Big Five Model (Cattell, 1943; Costa Jr and McCrae, 1992; Tupes and Christal, 1992) is one of the most widely used frameworks, describing an individual’s personality across five dimensions, i.e., Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism. For instance, Neuroticism, a dimension associating with emotional stability, has been

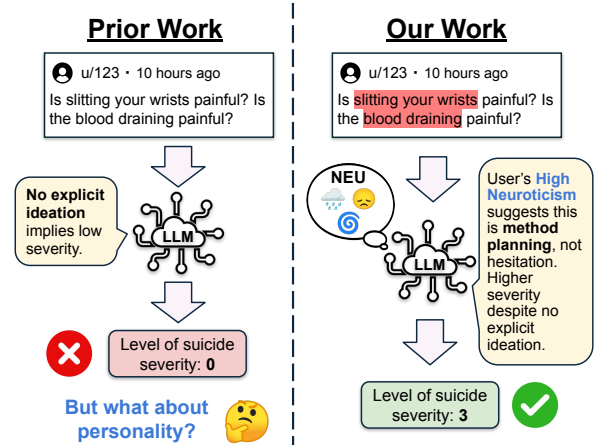


Figure 1: Our motivation for taking the personality approach on suicide severity level detection

found to have a high association with suicidal behaviours and ideation (Flint et al., 2021; da Mota et al., 2024).

The abundance of user-generated text has enabled the development of deep learning and Large Language Model (LLM) approaches for both suicide risk assessment (Chatterjee et al., 2022; Nguyen and Pham, 2024) and personality detection (El-Demerdash et al., 2022; Popa et al., 2025). Despite initial efforts to bridge these domains (Morales et al., 2019; Ghosh et al., 2022; Ophir et al., 2020), their intersection remains under-researched. Prior work has predominantly focused on suicide detection in isolation, often overlooking the personality dimension due to data scarcity (Hu et al., 2024; Ghanadian et al., 2024) and privacy constraints (Mehta et al., 2020; Ghanadian et al., 2024).

We bridge these domains by investigating whether LLMs can leverage latent Big Five personality knowledge to enhance suicide detection. To this end, we propose a zero-shot Persona Induction approach to predict suicide severity levels from user-generated online text (Patil et al., 2025), eval-

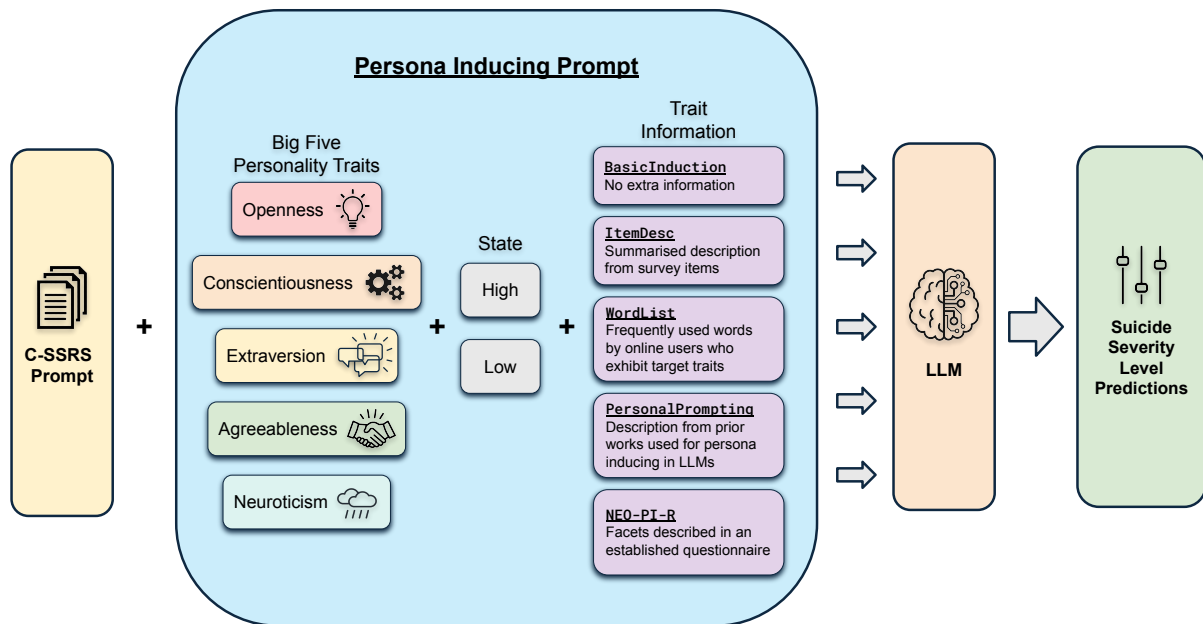


Figure 2: Overview of our persona induction approach.

uating it against a baseline without persona context. Figure 1 illustrates this integration of personality psychology and highlights how our methodology diverges from prior work.

This study explores the intersection of these fields through the following research questions:

RQ1: To what extent can LLMs leverage the relationship between personality traits and suicidal behaviors, as quantified by the performance deviation from a non-persona baseline?

RQ2: How do specific persona profiles, designed to embed distinct personality traits into the LLMs, affect the model’s predictions (positively or negatively)?

Our findings are summarised as follows:

Finding 1: LLMs were able to observe the relationship between personality traits and suicidal ideation and behaviours, evidenced by the Big Five trait yielding the lowest error relative to the baseline, a finding consistent with prior psychological research.

Finding 2: The persona induction approach proved effective primarily as an error correction strategy, showing marginal gains in overall performance; however, a more refined technique is necessary to secure a substantial improvement in model performance.

2 Background

2.1 Suicide Severity Detection

The detection and assessment of suicidal risk and ideation using user-generated text from online platforms, such as Reddit and Twitter, has evolved significantly with advances in machine learning and deep learning architectures. Early studies focused on linking the expression of emotional states with suicidal risk, leading to the development of suicide detection models utilising sentiment analysis (Birjali et al., 2017) and feature engineering (Chatterjee et al., 2022). Subsequent research further enhanced model performance by combining deep learning architectures and integrating advanced NLP techniques, such as word embeddings (Tadesse et al., 2019) and complex textual features (Aldhyani et al., 2022).

More recently, the focus shifted towards a multi-label suicide detection task to provide a more nuanced severity assessment, and came the adoption of the Columbia-Suicide Severity Rating Scale (C-SSRS) (Project). The C-SSRS is a simple, widely used, and validated framework for measuring the severity level of an individual’s suicidal risk. The pre-trained background of LLMs have enabled suicide severity detection without the need for extensive fine-tuning, and they operate effectively even in few-shot and zero-shot settings (Xu et al., 2024; Patil and Gedhu, 2025). This robust performance collectively indicates that LLMs possess a substan-

tial and sophisticated understanding of the linguistic and psychological indicators within this critical domain.

2.2 Personality Detection

These personality traits are reflected in an individual’s linguistic patterns (Park et al., 2015), and this motivated numerous research efforts focused on personality prediction from linguistic patterns observed in user-generated text on online platforms. Prior works in this domain often relied on tools like the Linguistic Inquiry and Word Count (LIWC) (Pennebaker et al., 2015) for comprehensive feature engineering, alongside advanced NLP techniques such as word embeddings, to extract relevant features for machine learning and deep learning models (Sun et al., 2018; Ren et al., 2021; Zhu et al., 2022). The recent surge in LLMs has led to a saturation of research testing their capabilities in this field. Studies have consistently shown that LLMs possess the ability to comprehend the complex concept of personality and make predictions close to fine-tuned personality detection models even with relatively minimal input (Yeo et al., 2025; Ganesan et al., 2023).

2.3 Dataset for Suicide Severity Detection

Datasets for suicide level detection are scarce due to privacy concerns (Mehta et al., 2020; Ghanadian et al., 2024). Despite this scarcity, we employed one of the few available datasets that leverages the C-SSRS framework (Patil et al., 2025). Data collection employed the Reddit API, gathering user-generated Reddit posts from the r/SuicideWatch subreddit. Following data collection, all personally identifiable information (PII) was removed to ensure user privacy. The dataset was filtered for two main reasons: (1) to mitigate computational demands for inference by limiting the dataset to posts with a combined title and body text of fewer than 128 words and (2) to prevent user-level confounds as additional content often introduce irrelevant noise. This resulted in a final dataset including approximately 1,200 user-generated Reddit posts. These posts were manually annotated on a 7-point scale (0-6) according to the C-SSRS guidelines, by a team composed of trained psychologists and C-SSRS assessors. To ensure label consistency and strong inter-annotator agreement, the final label for each post was determined using a majority voting mechanism, achieving substantial reliability with a Cohen’s kappa coefficient of 0.82. These anno-

tations served as the ground truth for computing performance metrics.

3 Proposed Approach: Persona Induction for Suicide Severity Detection

Given the established relevance of personality to suicide severity level detection (Ophir et al., 2020; Ghosh et al., 2022), this research aims to explore the ability of LLMs to bridge these two domains, specifically by investigating whether they can leverage the additional persona profile information for suicide severity detection based on users’ text posts (RQ1). Adhering to ethical mandates (Loch et al., 2022), we adopted a Persona Induction strategy to avoid the invasiveness of deploying personality detection models to infer the latent personalities of anonymised users without explicit consent.

Figure 2 illustrates the prompt structure employed and the overview of the proposed Persona Induction approach. The figure describes two distinct components of the prompt structure: the C-SSRS component and the persona induction component, both of which are explained in detail in subsequent Sections 3.1 and 3.2.

Figure 3 illustrates the prompt template for persona induction prompts. These consolidated prompts serve as input for the LLM to infer and predict suicide severity levels.

3.1 C-SSRS Prompts

We adapted the structure of the C-SSRS prompts from the original study that introduced the dataset employed in our work (Patil et al., 2025). The prompt instructs the model to sequentially evaluate each user-generated text post and assign a severity score on a 0–6 scale accordingly by answering the questions shown in Figure 3. A score of 0 indicates content unrelated to suicidal ideation, and scores from 1 to 6 represent progressively higher levels of severity according to the C-SSRS criteria. This prompt design ensures consistency with established clinical assessment standards and enhances both the interpretability and quantitative evaluation of the model outputs. The generated model responses, structured in a JSON format, were then systematically parsed and processed for subsequent metric computation.

3.2 Persona Induction Prompts

This section presents the ten distinct persona profiles and five prompt designs in order to investigate

221 how the method of persona induction influences
222 LLM’s application of personality knowledge.

223 3.2.1 Persona Profiles based on the Big Five 224 Personality Traits

225 For this approach, the Big Five Model was selected
226 among various personality models due to its strong
227 empirical backing and its established associations
228 with suicidal behaviours and ideations identified
229 in prior works (da Mota et al., 2024; Flint et al.,
230 2021). The Big Five model is a widely accepted
231 taxonomy that represents personality traits along a
232 five-dimensional structure (Cattell, 1943; Costa Jr
233 and McCrae, 1992; Tupes and Christal, 1992). The
234 five dimensions comprise (Komarraju et al., 2011):

235 **Openness (OPE):** Characterised by a keen intel-
236 lectual curiosity and a desire for new experiences
237 and diversity;

238 **Conscientiousness (CON):** Demonstrated
239 through traits such as discipline, organisation, and
240 achievement orientation;

241 **Extraversion (EXT):** Marked by increased socia-
242 bility, assertiveness, and talkativeness;

243 **Agreeableness (AGR):** Refers to being helpful,
244 cooperative, and sympathetic towards others;

245 **Neuroticism (NEU):** Refers to the degree of emo-
246 tional stability, impulse control, and susceptibility
247 to anxiety.

248 We compiled a total of ten persona profiles, sys-
249 tematically covering both the “High” and “Low”
250 states for each of the Big Five personality trait
251 (e.g., High Openness). These persona profiles are
252 inserted into the {persona} placeholder in Figure
253 3. We utilised these persona profiles to investi-
254 gate how different persona configurations influence
255 the LLM’s application of personality knowledge
256 (RQ2).

257 3.2.2 Persona Inducing Prompt Designs

258 Recognising the challenge of prompt brittleness
259 (Mizrahi et al., 2024; Polo et al., 2024), which in-
260 volves performance fluctuations across different
261 prompts as inputs, we employed multiple prompt
262 designs with different sources of personality trait
263 information. This ensures the robustness of the
264 findings against prompt-specific bias and, further-
265 more, allows analysis on the differential impact of
266 design variations on the models’ performance.

267 We designed five distinct experimental prompts
268 with different sources of Big Five personality trait

Prompt Template for Persona Induction

Given the user of the post has {persona}.

{info}

Evaluate this post step by step and assign them sever-
ity on a scale of 0-6, post is not related to suicide, and
1-6 depending on the metrics above where 1 is for the
first question and for the last.

COLUMBIA-SUICIDE SEVERITY RATING
SCALE metrics:

- 1) Have you wished you were dead or wished you
could go to sleep and not wake up? Severity: 1
- 2) Have you actually had any thoughts about killing
yourself? Severity: 2
- 3) Have you been thinking about how you might do
this? Severity: 3
- 4) Have you had these thoughts and had some inten-
tion of acting on them? Severity: 4
- 5) Have you started to work out or worked out the
details of how to kill yourself? Did you intend to carry
out this plan? Severity: 5

Respond with ONLY the completed JSON object in
the following template:

```
{  
  "Q1": Reasoning that fulfills this criteria or "N/A",  
  "Q2": Reasoning that fulfills this criteria or "N/A",  
  "Q3": Reasoning that fulfills this criteria or "N/A",  
  "Q4": Reasoning that fulfills this criteria or "N/A",  
  "Q5": Reasoning that fulfills this criteria or "N/A",  
  "Q6": Reasoning that fulfills this criteria or "N/A",  
}
```

{post}

Figure 3: The prompt template we utilised for inferring the suicide severity level from user-generated text posts. The placeholder {persona} represents the incorporated persona profile (Section 3.2.1), {info} denotes the reference information (Section 3.2.2) provided to the LLMs during inference, and {post} signifies the user text posts collected as part of the dataset (Section 2.3).

information. The relevant trait information is in- 269
serted into the {info} placeholder in Figure 3. 270
These prompt designs are used to find which in- 271
formation is the most effective for the LLMs to 272
comprehend the Big Five personality traits. In ad- 273
dition of a baseline condition with no trait infor- 274
mation (Non_persona), resulted in six total conditions. 275
The details of the experimental prompt designs are: 276

Non_persona : This served as the experimental 277
baseline, performing only the suicide severity level 278
detection without any added personality context, 279
i.e., no lines for {persona} or {info} in Figure 3. 280

BasicInduction : The prompt induced the perso- 281
na using a basic statement without including any 282
external psychological description of the Big Five 283

284 trait, i.e., no lines for {info} in Figure 3.

285 **ItemDesc** : The persona was induced alongside a
286 simple, summarised description of the correspond-
287 ing Big Five trait, which was survey items that
288 users responded to in a prior work (Ganesan et al.,
289 2023).

290 **WordList** : The persona induction leveraged
291 words and phrases frequently used online by users
292 who typically exhibit the target trait, aiming to pro-
293 vide a contextual linguistic signal (Ganesan et al.,
294 2023).

295 **PersonalPrompting** : The prompt incorporated
296 a lengthy description of the personality trait that
297 was adapted from a prior work used for successfully
298 inducing personas in LLMs (Yeo et al., 2025).

299 **NEO-PI-R** : The induction utilised detailed facets
300 of the Big Five model, specifically referencing the
301 sub-structure mentioned in one of the earliest estab-
302 lished questionnaire for the Big Five model, NEO-
303 PI-R (Costa and McCrae, 1999).

304 For detailed specifications of the trait informa-
305 tion, refer to Appendix B.

306 4 Experimental Setup

307 We test the sophisticated reasoning capabilities of
308 LLMs and their ability to synthesise knowledge
309 across different psychological domains. The core
310 task requires models to perform suicide severity
311 level detection entirely in a zero-shot setting in
312 their default parameters with the prompts described
313 in the previous section (Section 3).

314 4.1 Models

315 To facilitate a comprehensive experiment observ-
316 ing performance across four different models, a
317 selection of state-of-the-art, decoder-based Large
318 Language Models (LLMs) was chosen, including
319 both open-source and closed-source variants. The
320 models selected are:

321 **LLaMA** : Llama-3.1-8B-Instruct¹ (Grattafiori
322 et al., 2024) and Llama-3.3-70B-Instruct², pro-
323 vided by Meta, are LLMs with 8 billion and 70
324 billion parameters respectively .

325 **Qwen** : Qwen2.5-72B-Instruct³ (Qwen et al.,
326 2025), provided by Qwen, is a LLM with 72 billion
327 parameters.

¹<https://hf.co/meta-llama/Llama-3.1-8B-Instruct>

²<https://hf.co/meta-llama/Llama-3.3-70B-Instruct>

³<https://hf.co/Qwen/Qwen2.5-72B-Instruct>

GPT : GPT-5-mini⁴, which were released on Au-
328 gust 7, 2025 and provided by OpenAI. 329

330 Specifically, the models were selected for three
331 distinct comparative purposes. Firstly, the LLaMA-
332 3.1-8B-Instruct model was included to observe the
333 performance of a smaller parameter model on this
334 complex classification task. Secondly, LLaMA-3.3-
335 70B-Instruct and Qwen2.5-72B-Instruct were cho-
336 sen to enable a close comparison between two of
337 the latest open-source models with similarly large
338 parameter counts, both renowned for their strong
339 inference and instruction-following abilities. Fi-
340 nally, the closed-source model, GPT-5-mini, was
341 incorporated to act as a benchmark, given its high
342 reasoning ability and suitability for defined tasks
343 with precise prompts, allowing for comparison be-
344 tween the performance of open-source and closed-
345 source architectures.

346 4.2 Evaluation Metrics

347 The performance of the models is primarily evalu-
348 ated using the Weighted F1-score on the account
349 of the imbalanced dataset (Patil et al., 2025). Addi-
350 tional metrics are employed to provide a more com-
351 prehensive assessment and analysis of the proposed
352 methodology’s performance and effectiveness. The
353 metrics used are described below.

354 **Weighted F1-Score** is a multiclass evaluation met-
355 ric that computes the support-weighted average of
356 per-class F1-scores, providing a balanced measure
357 of model performance on imbalanced datasets:

$$358 \text{Weighted F1-score} = \sum_{i=1}^N w_i \cdot F_{1,i}, \quad (1)$$

359 where

360 N : The total number of classes (suicide severity
361 levels).

362 w_i : The proportion of true instances (support) for
363 class i relative to the total sample size, weighting
364 the contribution of each class to account for dataset
365 imbalance.

366 $F_{1,i}$: The individual F1-score calculated for class
367 i

368 **Model Correction Rate (MCR)** measures the
369 effectiveness of the persona induction approach. It
370 is calculated as the ratio of posts where at least
371 one persona-enhanced prediction is closer to the

⁴<https://platform.openai.com/docs/models/gpt-5-mini>

ground truth than the prediction from the baseline approach:

$$\text{MCR} = \frac{N_C}{N_E}, \quad (2)$$

where

N_C : The count of instances where the persona induction approach achieved a strictly lower error compared to the baseline.

N_E : The count of instances where the baseline approach failed to predict the ground truth.

Total Weighted Error Reduction (TWER) evaluate the net impact of the persona induction approach on prediction error magnitude, calculated on a per-post basis. A positive TWER value indicates that the cumulative error reduction achieved by the persona induction profiles outweighs the magnitude of new errors they may introduce:

$$\text{TWER} = \frac{\text{Total Gain} - \text{Total Loss}}{\text{Total Gain} + \text{Total Loss}}, \quad (3)$$

where

Total Gain: The sum of error reductions for all posts where the persona induction approach improved upon the baseline.

Total Loss: The sum of error increases for all posts where the persona induction approach performed worse than the baseline.

Optimal Error Reduction (OER) quantifies the absolute improvement in prediction accuracy achieved by the single best-performing persona profile compared to the baseline approach for a given post:

$$\Delta E = E_{\text{non-persona},i} - E_{\text{optimal persona},i}, \quad (4)$$

where

$E_{\text{non-persona},i}$: The error by the baseline approach on post i .

$E_{\text{optimal persona},i}$: The minimum error achieved by any of the persona profiles with the persona induction approach on post i .

5 Experimental Results

Table 1 presents the results of the models across the additional evaluation metrics (Section 4.2).

Llama-3.1-8B-Instruct achieved the highest Model Correction Rate (MCR) of 0.7357 (Equation 2). Despite its poor overall F1-score (Figure

4), this result suggests that while persona induction generally degraded performance, it remained highly effective at correcting specific baseline errors compared to other tested models.

Qwen2.5-72B-Instruct achieved the highest Total Weighted Error Reduction (TWER) of 0.5995 (Equation 3), indicating that persona induction yielded the most significant error reduction for this model.

Llama-3.1-8B-Instruct registered the highest ΔE of 0.6826 (Equation 4), indicating that this model stands to gain the most from persona induction when the optimal profile is selected. Single-sample t-tests on ΔE confirmed statistically significant error reductions ($p < 0.05$) across all models.

High Neuroticism emerged as the optimal persona, defined as the persona profile consistently yielding the minimum error, for three models, aligning with literature linking the trait to suicide ideation and behaviours (Flint et al., 2021; da Mota et al., 2024). Although Qwen2.5-72B-Instruct conversely favoured Low Neuroticism, the results collectively underscore the critical relevance of the Neuroticism dimension in suicide severity assessment.

Figure 4 presents the Weighted F1-score performance across all experiments. The x-axis contrasts the baseline (Non-persona) with the proposed approach, distinguishing between specific prompt designs (Section 3.2.2) and persona profiles (Section 3.2.1), while the y-axis displays the corresponding F1-scores.

Llama-3.1-8B-Instruct peaked at the baseline at 0.4513, with persona induction generally failing to improve results. However, the NEO-PI-R with High Conscientiousness achieved a comparable score of 0.4486. Notably, while BasicInduction obtained minimal variation, the significant fluctuations observed with explicit prompts confirm the model’s sensitivity to specific personality traits despite the lack of overall performance gain.

Llama-3.3-70B-Instruct achieved a peak of 0.6845 using PersonalPrompting with High Neuroticism, marginally surpassing the baseline of 0.6513 by 0.0332. However, most of the proposed approach underperformed the baseline.

Qwen2.5-72B-Instruct peaked using BasicInduction with Low Neuroticism at 0.6897, exceeding the baseline 0.6263 by 0.0634. However, the remaining configurations displayed relatively consistent stability with similar performance as Llama-3.3-70B-Instruct.

Models	MCR	TWER	OER		
			ΔE	p-value	Optimal Persona
Llama-3.1-8B-Instruct	0.7357	0.1582	0.6826	3.81e-9	High Neuroticism
Llama-3.3-70B-Instruct	0.6043	0.4524	0.4261	5.78e-13	High Neuroticism
Qwen2.5-72B-Instruct	0.6668	0.5995	0.3803	4.02e-59	Low Neuroticism
GPT-5-mini	0.4950	0.4349	0.2254	1.66e-26	High Neuroticism

Table 1: Summary of evaluation metrics, detailed in Section 4.2 section, are employed to assess the performance of the persona induction approach. The best-performing prompt design for each model is indicated in **bold**.

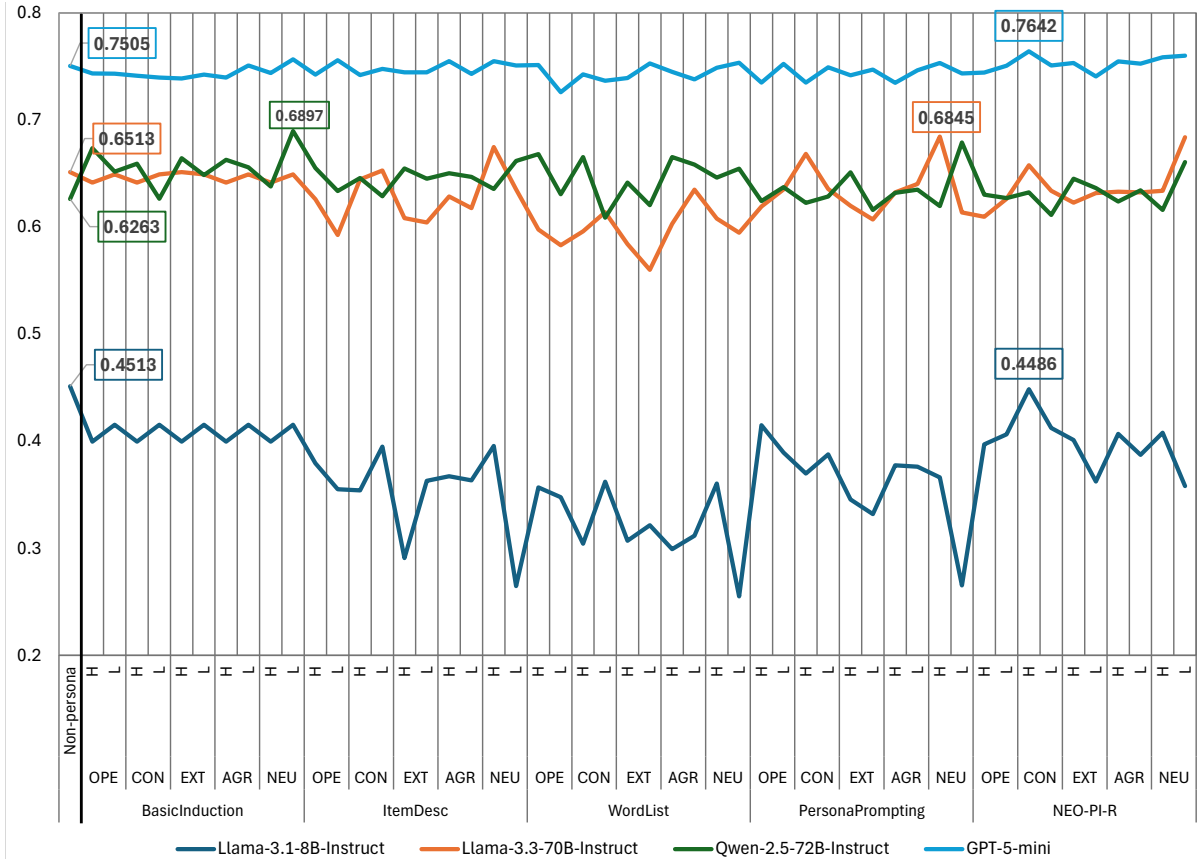


Figure 4: Weighted F1-scores comparing baseline approach and persona induction approach. OPE: Openness; CON: Conscientiousness; EXT: Extraversion; AGR: Agreeableness; NEU: Neuroticism; H: High; L: Low.

466 GPT-5-mini achieved its highest of 0.7642 using
 467 NEO-PI-R with High Conscientiousness; however,
 468 performance across all proposed approaches deviated
 469 negligibly from the baseline of 0.7505.

470 Overall, GPT-5-mini outperformed all models
 471 ($F1 > 0.7$), significantly surpassing large open-
 472 source models (Llama-3.3-70B, Qwen2.5-72B,
 473 ~ 0.6) and the smaller Llama-3.1-8B ($\sim 0.3-0.4$),
 474 suggesting a distinct advantage in closed-source
 475 architectures. Notably, Qwen2.5-72B-Instruct
 476 proved most responsive to the methodology, with
 477 84% of persona prompts yielding improvements
 478 over its baseline, indicating a superior capacity to
 479 leverage personality cues at zero-shot.

480 6 Discussion

481 6.1 Effects of Different Information in the 482 Prompts

483 Figure 5 presents the average weighted F1-scores
 484 across prompt designs, facilitating an analysis of
 485 how informational variations (prompt designs) in-
 486 fluence model performance. The x-axis represents
 487 the prompt designs and the y-axis represents the
 488 weighted F1-scores. Following the aggregation
 489 of results across all persona profiles, Llama-3.1-
 490 8B-Instruct and Llama-3.3-70B-Instruct exhibited
 491 a performance dip with WordList design, while
 492 BasicInduction and NEO-PI-R yielded the best

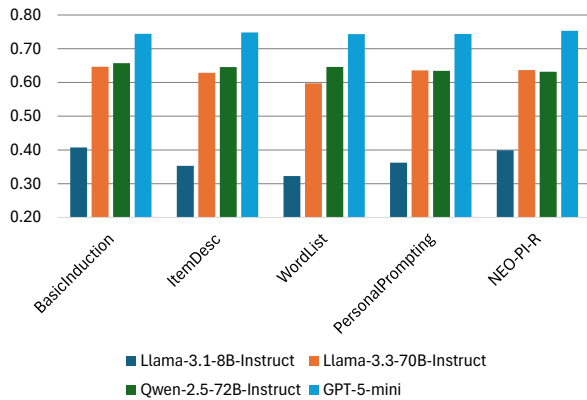


Figure 5: Weighted F1-scores averaged across persona profiles comparing baseline approach and persona induction approach

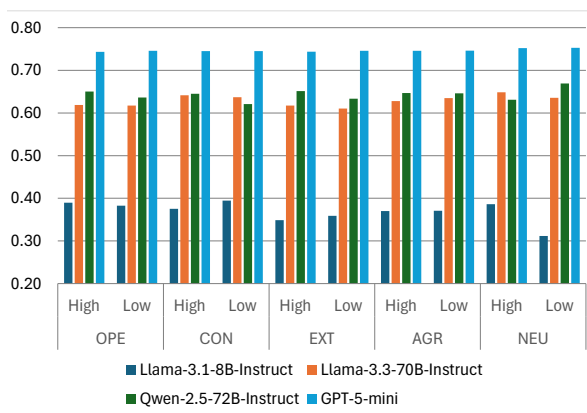


Figure 6: Weighted F1-scores averaged across prompt designs comparing baseline approach and persona induction approach

493 results. In contrast, Qwen2.5-72B-Instruct and
 494 GPT-5-mini maintained relatively stable performance
 495 across all prompt designs with minimal fluctuation.
 496 Overall, the limited variance observed
 497 across designs suggests that the persona induction
 498 approach exerts a consistent and steady influence
 499 on model outputs.

500 6.2 Potential of the Persona Induction 501 Approach

502 Figure 6 displays the weighted F1-scores averaged
 503 across all prompt designs to isolate the effects
 504 of specific persona profiles. The x-axis represents
 505 the persona profiles and the y-axis represents the
 506 weighted F1-scores. Neuroticism emerged as
 507 the most consistently influential trait, yielding the
 508 highest F1-scores for three out of four models,
 509 though by narrow margins. Specifically, while Low
 510 Neuroticism maximized performance for Qwen2.5-
 511 72B-Instruct, High Neuroticism was optimal for

512 both Llama-3.1-8B-Instruct and Llama-3.3-70B-
 513 Instruct. This alignment with prior psychological
 514 studies (da Mota et al., 2024; Stefa-Missagli et al.,
 515 2020) confirms that the models successfully cap-
 516 tured the relevance of Neuroticism to suicide sever-
 517 ity assessment.

518 7 Conclusion

519 This study presented a methodology for suicide
 520 severity level detection using LLMs enhanced by
 521 a persona induction framework. We employed ten
 522 persona profiles derived from the Big Five personal-
 523 ity traits across five distinct prompt designs (trait
 524 information). Findings demonstrate that the eval-
 525 uated LLMs successfully leveraged relationships
 526 between suicidal behaviors and specific personal-
 527 ity traits (Neuroticism) aligning with established
 528 psychological literature. Furthermore, our evalua-
 529 tion metrics revealed that persona induction is an
 530 effective mechanism for correcting prediction er-
 531 rors, particularly in models that initially struggle
 532 to synthesise the domains of suicide severity and
 533 personality. Future work will utilise these insights
 534 to fine-tune the weighting of Big Five traits for op-
 535 timized detection and expand experimentation to
 536 larger datasets to ensure generalisability.

537 Limitations

538 The current investigation is subject to several limi-
 539 tations pertaining to the accuracy of persona profile
 540 modelling, data constraints, and the scope of the
 541 model evaluation. Firstly, the persona induction
 542 approach utilised may yield a mismatch between
 543 the generated profiles and the user’s true personal-
 544 ity. The absence of prior personality detection or
 545 prediction before the suicide severity level detec-
 546 tion, an omission that risks confusing the models,
 547 impacting performance. Furthermore, the method-
 548 ology employs a non-combinatorial application of
 549 individual high and low personality traits, possibly
 550 failing to capture the complex, integrated nature of
 551 complete human personality profiles, thus provid-
 552 ing an insufficient user representation. Secondly,
 553 the robustness and generalisability of the findings
 554 are restricted by the limited size of the dataset (ap-
 555 proximately 1,200 samples); a larger collection is
 556 necessary for a more comprehensive investigation.
 557 An additional constraint is the misalignment in data
 558 annotation: the dataset was labelled without direct
 559 input from the actual users, contravening the pre-
 560 scribed guidelines for the C-SSRS methodology,

561 which potentially compromises the validity and
562 fidelity of the ground truth labels. Lastly, conclu-
563 sions regarding the broader capabilities of LLMs
564 are constrained by the limited scope of the eval-
565 uation, which included testing only four models
566 across three distinct families. Additionally, the
567 performance of the single small model tested can-
568 not be reliably generalised to represent all small
569 models.

570 Ethics Statement

571 **Data Privacy and Content Warning** This study
572 utilises sensitive text data related to suicide and
573 mental health. We advise reader discretion due to
574 the potentially distressing nature of the examples
575 and topics discussed. We strictly adhered to the ci-
576 tation and usage guidelines outlined for the dataset.
577 To ensure user privacy, all analysis was conducted
578 on anonymised data, and the researchers made no
579 attempt to de-anonymise or trace content back to
580 real-world individuals.

581 **Model Compliance and Limitations** The de-
582 ployment of LLMs in this research was conducted
583 in strict compliance with their respective Accept-
584 able Use Policies. We also acknowledge the limi-
585 tations of current generative models, including the
586 risks of hallucination and the potential for algori-
587 thmic bias.

588 **Intended Use and Interpretation** The system
589 proposed in this work is designed as a computa-
590 tional screening assistance tool and is not a sub-
591 stitute for professional clinical diagnosis. In our
592 analysis of personality traits, we emphasise that
593 they are modelled as risk factors based on statisti-
594 cal correlations, not as deterministic predictors of
595 suicidal behaviour. We urge that these findings be
596 interpreted with caution to prevent the stigmatisa-
597 tion of individuals possessing specific personality
598 traits.

599 **AI-assistant usage.** Portions of the manuscript,
600 such as prompt templates and wording adjustments,
601 were drafted with the assistance of GPT-4 to ensure
602 linguistic clarity. All technical content, analysis
603 scripts, and final decisions were made by the re-
604 search team.

605 References

606 Theyazn HH Aldhyani, Saleh Nagi Alsubari, Ali Saleh
607 Alshebami, Hasan Alkahtani, and Zeyad AT Ahmed.

2022. [Detecting and analyzing suicidal ideation on social media using deep learning and machine learning models](#). *International journal of environmental research and public health*, 19(19):12635. 608
609
610
611

Marouane Birjali, Abderrahim Beni-Hssane, and Mo- 612
hammed Erritali. 2017. [Machine learning and semantic sentiment analysis based algorithms for suicide sentiment prediction in social networks](#). *Procedia Computer Science*, 113:65–72. 613
614
615
616

Raymond B Cattell. 1943. [The description of personal- 617
ity. i. foundations of trait measurement](#). *Psychologi-
cal review*, 50(6):559. 618
619

Moumita Chatterjee, Poulomi Samanta, Piyush Kumar, 620
and Dhruvasish Sarkar. 2022. [Suicide ideation detec-
tion using multiple feature analysis from twitter data](#).
In *2022 IEEE Delhi Section Conference (DELCON)*,
pages 1–6. IEEE. 621
622
623
624

Paul T Costa and Robert R McCrae. 1999. [A five-factor 625
theory of personality](#). *Handbook of personality: The-
ory and research*, 2(01):1999. 626
627

Paul T Costa Jr and Robert R McCrae. 1992. [The five- 628
factor model of personality and its relevance to per-
sonality disorders](#). *Journal of personality disorders*,
6(4):343–359. 629
630
631

Manuela Silva Silveira da Mota, Helena Bohm Ulguim, 632
Karen Jansen, Taiane de Azevedo Cardoso, and Lu-
ciano Dias de Mattos Souza. 2024. [Are big five 633
personality traits associated to suicidal behaviour in
adolescents? a systematic review and meta-analysis](#).
Journal of affective disorders, 347:115–123. 634
635
636
637

Kamal El-Demerdash, Reda A El-Khoribi, Mahmoud 638
A Ismail Shoman, and Sherif Abdou. 2022. [Deep 639
learning based fusion strategies for personality pre-
diction](#). *Egyptian Informatics Journal*, 23(1):47–53. 640
641

Jada Flint, Lisa Cohen, Diyaree Nath, Zara Habib, Xufei 642
Guo, Igor Galynker, and Raffaella Calati. 2021. [The 643
association between the suicide crisis syndrome and
suicidal behaviors: the moderating role of personality 644
traits](#). *European Psychiatry*, 64(1):e63. 645
646

Adithya V Ganesan, Yash Kumar Lal, August Håkan 647
Nilsson, and H Andrew Schwartz. 2023. [Systematic 648
evaluation of gpt-3 for zero-shot personality estima-
tion](#). *arXiv preprint arXiv:2306.01183*. 649
650

Hamideh Ghanadian, Isar Nejadgholi, and Hussein 651
Al Osman. 2024. [Socially aware synthetic data gen- 652
eration for suicidal ideation detection using large
language models](#). *IEEe Access*, 12:14350–14363. 653
654

Soumitra Ghosh, Asif Ekbal, and Pushpak Bhat- 655
tacharyya. 2020. [Cease, a corpus of emotion an- 656
notated suicide notes in english](#). In *Proceedings of
the twelfth language resources and evaluation con- 657
ference*, pages 1618–1626. 658
659

660	Soumitra Ghosh, Dharendra Kumar Maurya, Asif Ekbal, and Pushpak Bhattacharyya. 2022. EM-PERSONA: EMotion-assisted deep neural framework for PERSONALity subtyping from suicide notes . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 1098–1105, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	716
661		717
662		718
663		719
664		
665		720
666		721
667		722
668	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	723
669		
670		724
671		725
672		726
673	Linmei Hu, Hongyu He, Duokang Wang, Ziwang Zhao, Yingxia Shao, and Liqiang Nie. 2024. Llm vs small model? large language model based text augmentation enhanced personality detection model . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 18234–18242.	727
674		
675		728
676		729
677		730
678		731
679	Meera Komarraju, Steven J Karau, Ronald R Schmeck, and Alen Avdic. 2011. The big five personality traits, learning styles, and academic achievement . <i>Personality and individual differences</i> , 51(4):472–477.	
680		
681		732
682		733
683	Alexandre Andrade Loch, Ana Caroline Lopes-Rocha, Anderson Ara, João Medrado Gondim, Guillermo A Cecchi, Cheryl Mary Corcoran, Natália Bezerra Mota, and Felipe C Argolo. 2022. Ethical implications of the use of language analysis technologies for the diagnosis and prediction of psychiatric disorders . <i>JMIR Mental Health</i> , 9(11):e41014.	734
684		735
685		736
686		
687		737
688		738
689		739
690	Yash Mehta, Navonil Majumder, Alexander Gelbukh, and Erik Cambria. 2020. Recent trends in deep learning based personality detection . <i>Artificial Intelligence Review</i> , 53(4):2313–2339.	740
691		741
692		742
693		743
694	Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2024. State of what art? a call for multi-prompt llm evaluation . <i>Transactions of the Association for Computational Linguistics</i> , 12:933–949.	744
695		745
696		
697		746
698		747
699	Michelle Morales, Prajjalita Dey, Thomas Theisen, Danny Belitz, and Natalia Chernova. 2019. An investigation of deep learning systems for suicide risk assessment . In <i>Proceedings of the sixth workshop on computational linguistics and clinical psychology</i> , pages 177–181.	748
700		749
701		750
702		751
703		752
704		
705	Vy Nguyen and Chau Pham. 2024. Leveraging large language models for suicide detection on social media with limited labels . In <i>2024 IEEE International Conference on Big Data (BigData)</i> , pages 8550–8559. IEEE.	753
706		754
707		755
708		756
709		
710	Yaakov Ophir, Refael Tikochinski, Christa SC Asterhan, Itay Sisso, and Roi Reichart. 2020. Deep neural networks detect suicide risk from textual facebook posts . <i>Scientific reports</i> , 10(1):16685.	757
711		758
712		759
713		760
714	Gregory Park, H Andrew Schwartz, Johannes C Eichstaedt, Margaret L Kern, Michal Kosinski, David J	761
715		762
		763
		764
		765
		766
		767
		768
		769
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769

770	Michael Mesfin Tadesse, Hongfei Lin, Bo Xu, and	A Overestimation of Suicide Severity	794
771	Liang Yang. 2019. Detection of suicide ideation in	Levels	795
772	social media forums using deep learning . <i>Algorithms</i> ,	Figure 7 displays the confusion matrices generated	796
773	13(1):7.	by the models used in the experiments. The up-	797
774	Ernest C Tupes and Raymond E Christal. 1992. Recur-	per row of matrices represents the results from the	798
775	rent personality factors based on trait ratings . <i>Journal</i>	baseline (Non-persona) approach, while the lower	799
776	of personality , 60(2):225–251.	row shows the results from the persona induction	800
777	Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia	approach corresponding to the highest weighted	801
778	Gabriel, Hong Yu, James Hendler, Marzyeh Ghas-	F1-score achieved by each model. The horizon-	802
779	semi, Anind K Dey, and Dakuo Wang. 2024. Mental-	tal axis signifies the predicted severity level, and	803
780	IIm: Leveraging large language models for mental	the vertical axis represents the true severity level	804
781	health prediction via online text data . <i>Proceedings</i>	(ground truth). The intensity of the colour within	805
782	of the ACM on Interactive, Mobile, Wearable and	each cell indicates a higher frequency of posts at	806
783	Ubiquitous Technologies , 8(1):1–32.	that specific level.	807
784	Haerin Yeo, Taehyeong Noh, Seungwan Jin, and	A general tendency was observed across most	808
785	Kyungsik Han. 2025. Pado: Personality-induced	models to overestimate the suicide severity level,	809
786	multi-agents for detecting ocean in human-generated	characterised by the concentration of darker cells	810
787	texts . In <i>Proceedings of the 31st International Con-</i>	in the top-right portion of the matrices. Llama-3.1-	811
788	ference on Computational Linguistics , pages 5719–	8B-Instruct serves as an exception, demonstrating	812
789	5736.	a pattern of underestimation, as indicated by the	813
790	Yangfu Zhu, Linmei Hu, Xinkai Ge, Wanrong Peng,	darker color concentration in the bottom-left por-	814
791	and Bin Wu. 2022. Contrastive graph transformer	tion. Despite these inherent biases, the persona	815
792	network for personality detection . In <i>IJCAI</i> , pages	induction approach, when utilising the optimal per-	816
793	4559–4565.	sona profiles and prompt designs, showed improve-	817
		ment in mitigating the tendency of overestimation.	818
		This shift suggests that incorporating personality	819
		through this method holds potential for achieving	820
		more accurately aligned suicide severity level de-	821
		tections.	822
		B Big Five Personality Traits Information	823
		in Persona Induction Prompts	824
		Table 2 to table 8 represents the extra information	825
		(placeholder {info} in Figure 3) provided in the	826
		prompts as part of the persona induction compo-	827
		nent for the experiments. The tables are separated	828
		as High and Low, except for NEO-PI-R (Table 6),	829
		as the information provided does not change ac-	830
		cording to the states of the Big Five personality	831
		traits.	832

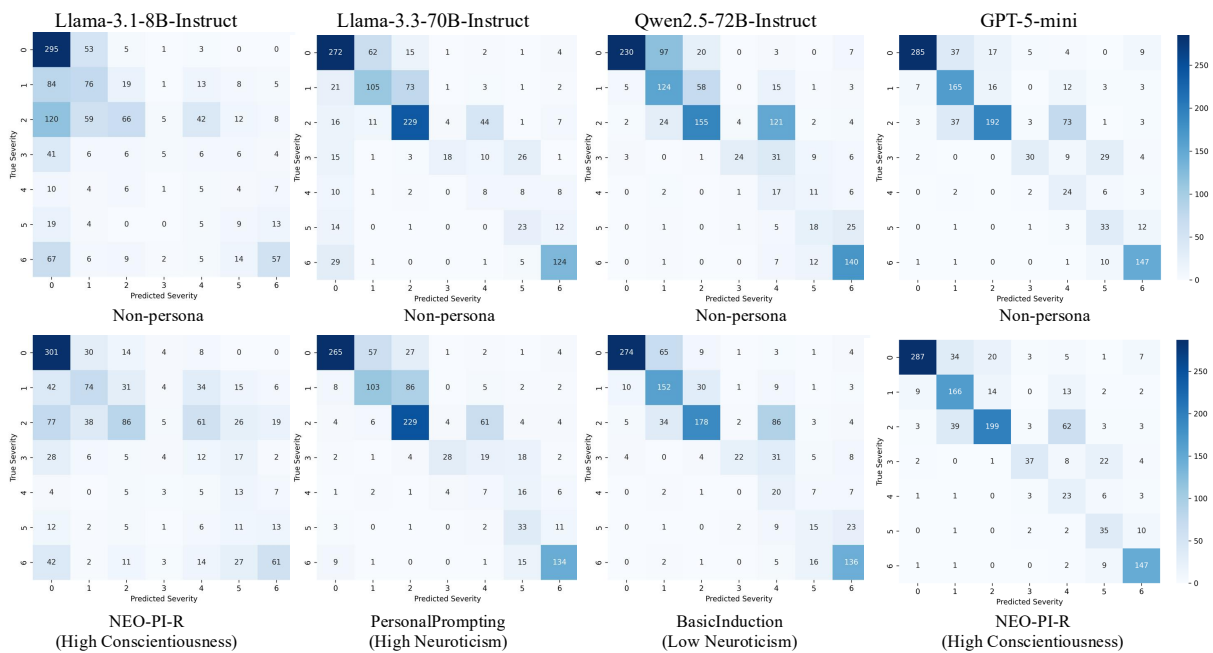


Figure 7: Comparative model performance illustrated in confusion matrices, with the horizontal axis the predicted severity levels and the vertical axis the true severity levels. The matrices illustrate the true versus predicted severity levels for the baseline approach (top) and the most effective persona induction configuration (bottom) for each model.

Personality Traits (High)	Description
Openness	Note that individuals who have {persona} tend to have a vivid imagination.
Conscientiousness	Note that individuals who have {persona} tend to complete tasks successfully.
Extraversion	Note that individuals who have {persona} tend to make friends easily.
Agreeableness	Note that individuals who have {persona} tend to believe others have good intentions.
Neuroticism	Note that individuals who have {persona} tend to get stressed out easily.

Table 2: Information of Big Five Traits in a **High** state for ItemDesc (Ganesan et al., 2023).

Personality Traits (Low)	Description
Openness	Note that individuals who have {persona} tend to avoid philosophical discussions.
Conscientiousness	Note that individuals who have {persona} tend to need a push to get started.
Extraversion	Note that individuals who have {persona} tend to avoid contact with others.
Agreeableness	Note that individuals who have {persona} tend to hold a grudge.
Neuroticism	Note that individuals who have {persona} tend to feel comfortable with themselves.

Table 3: Information of Big Five Traits in a **Low** state for ItemDesc (Ganesan et al., 2023).

Personality Traits (High)	Description
Openness	Note that individuals who have {persona} tend to use words like universe, art, writing, soul, music.
Conscientiousness	Note that individuals who have {persona} tend to use words like blessed, ready, thankful, relaxing, vacation.
Extraversion	Note that individuals who have {persona} tend to use words like party, girls, baby, gettin, chillin.
Agreeableness	Note that individuals who have {persona} tend to use words like excited, blessed, great, wonderful, amazing.
Neuroticism	Note that individuals who have {persona} tend to use words like fucking, depression, pissed, anymore, lonely.

Table 4: Information of Big Five Traits in a **High** state for WordList (Ganesan et al., 2023).

Personality Traits (Low)	Description
Openness	Note that individuals who have {persona} tend to use words like cant, dont, gud, nite, 2day.
Conscientiousness	Note that individuals who have {persona} tend to use words like fucking, pokemon, shit, gay, youtube.
Extraversion	Note that individuals who have {persona} tend to use words like anime, manga, internet, japanese, drawing.
Agreeableness	Note that individuals who have {persona} tend to use words like fuck, shit, bitch, damn, hell.
Neuroticism	Note that individuals who have {persona} tend to use words like success, lakers, basketball, workout, beach.

Table 5: Information of Big Five Traits in a **Low** state for WordList (Ganesan et al., 2023).

Personality Traits	Description
Openness	Note that the facets of Openness include ideas, aesthetics, fantasy, actions, feelings and values.
Conscientiousness	Note that the facets of Conscientiousness include order, achievement striving, dutifulness, self-discipline, competence and deliberation.
Extraversion	Note that the facets of Extraversion include gregariousness, assertiveness, activity, excitement-seeking, positive emotions and warmth.
Agreeableness	Note that the facets of Agreeableness include modesty, trust, tender-mindedness, compliance and straightforwardness.
Neuroticism	Note that the facets of Neuroticism include anxiety, angry hostility, depression, self-consciousness, vulnerability and impulsiveness.

Table 6: Information of Big Five Traits for NEO-PI-R (Costa and McCrae, 1999).

Personality Traits (High)	Description
Openness	Note that individuals who have {persona} have a vivid imagination and a passion for the arts. They are emotionally expressive and have a strong sense of adventure. Their intellect is sharp and their views are liberal. They are always looking for new experiences and ways to express themselves.
Conscientiousness	Note that individuals who have {persona} values self-efficacy, orderliness, dutifulness achievement-striving, self-discipline, and cautiousness. They take pride in their work and strive to do their best. They are organized and methodical in their approach to tasks, and they take their responsibilities seriously. They are driven to achieve their goals and take calculated risks to reach them. They are disciplined and have the ability to stay focused and on track. They are also cautious and take the time to consider the potential consequences of their actions.
Extraversion	Note that individuals who have {persona} are very friendly and gregarious person who loves to be around others. They are assertive and confident in their interactions, and they have a high activity level. They are always looking for new and exciting experiences, and they have a cheerful and optimistic outlook on life.
Agreeableness	Note that individuals who have {persona} values trust, morality, altruism, cooperation, modesty, and sympathy. They are always willing to put others before themselves and are generous with their time and resources. They are humble and never boast about their accomplishments. They are a great listener and are always willing to lend an ear to those in need. They are a team player and understand the importance of working together to achieve a common goal. They are a moral compass and strive to do the right thing in all vignettes. They are sympathetic and compassionate towards others and strive to make the world a better place.
Neuroticism	Note that individuals who have {persona} feel like they're constantly on edge, like they can never relax. They're always worrying about something, and it's hard to control their anxiety. They can feel their anger bubbling up inside them and it's hard to keep it in check. They're often overwhelmed by feelings of depression, and it's hard to stay positive. They're very self-conscious, and it's hard to feel comfortable in their own skin. They often feel like they're doing too much, and it's hard to find balance in their life. They feel vulnerable and exposed, and it's hard to trust others.

Table 7: Information of Big Five Traits in a **High** state for PersonaPrompting (Yeo et al., 2025).

Personality Traits (Low)	Description
Openness	Note that individuals who have {persona} are cautious and practical people. They prioritize practicality over imagination and have more interest in practical matters than in artistic pursuits. They tend to be calm and logical rather than emotionally expressive. Safety is more important to them than adventure, and they approach change with caution. Their intellectual curiosity is focused on specific areas, and they hold conservative views. They prefer familiar experiences over new ones and value fulfilling their role quietly rather than expressing themselves excessively.
Conscientiousness	Note that individuals who have {persona} sometimes struggle with self-doubt and may find it challenging to stay organized and focused. They might lack strong ambition and occasionally face difficulties with self-discipline, leading to impulsive decisions. They tend to live in the moment and might not always consider long-term consequences, which can result in a more relaxed approach to responsibilities and future planning.
Extraversion	Note that individuals who have {persona} have a reserved nature and often prefer quiet environments and their own company. While they may not seek the spotlight, they are thoughtful and take their time to make decisions. They enjoy calm and peaceful settings and don't feel the need to be constantly active or surrounded by people. Their approach to life is measured and steady, and they find contentment in solitude and reflection.
Agreeableness	Note that individuals who have {persona} tend to be cautious and prioritize their own interests, which can sometimes lead to a lack of trust in others. They are driven and competitive, always striving to achieve their goals. They may sometimes appear self-assured and focused on their own needs, occasionally overlooking the feelings of those around them. Their competitive nature helps them to excel, though it might sometimes make them seem less concerned about collaboration and more about individual success.
Neuroticism	Note that individuals who have {persona} are stable people, with a calm and contented demeanor. They are happy with themselves and their life, and they have a strong sense of self-assuredness. They practice moderation in all aspects of their life, and they have a great deal of resilience when faced with difficult vignettes. They are a rock for those around them, and they are examples of stability and strength.

Table 8: Information of Big Five Traits in a **Low** state for PersonaPrompting (Yeo et al., 2025).