

# EmoPerso: Enhancing Personality Detection with Self-Supervised Emotion-Aware Modelling

Lingzhi Shen

School of Electronics and Computer  
Science, University of Southampton  
Southampton, United Kingdom  
l.shen@soton.ac.uk

Xiaohao Cai

School of Electronics and Computer  
Science, University of Southampton  
Southampton, United Kingdom  
x.cai@soton.ac.uk

Yunfei Long

School of Electronic Engineering and  
Computer Science, Queen Mary  
University of London  
London, United Kingdom  
yunfei.long@qmul.ac.uk

Imran Razzak

Department of Computational  
Biology, Mohamed bin Zayed  
University of Artificial Intelligence  
Abu Dhabi, United Arab Emirates  
imran.razzak@mbzuai.ac.ae

Guanming Chen

School of Electronics and Computer  
Science, University of Southampton  
Southampton, United Kingdom  
gc3n21@soton.ac.uk

Shoaib Jameel

School of Electronics and Computer  
Science, University of Southampton  
Southampton, United Kingdom  
M.S.Jameel@southampton.ac.uk

## Abstract

Personality detection from text is commonly performed by analysing users' social media posts. However, existing methods heavily rely on large-scale annotated datasets, making it challenging to obtain high-quality personality labels. Moreover, most studies treat emotion and personality as independent variables, overlooking their interactions. In this paper, we propose a novel self-supervised framework, EmoPerso, which improves personality detection through emotion-aware modelling. EmoPerso first leverages generative mechanisms for synthetic data augmentation and rich representation learning. It then extracts pseudo-labeled emotion features and jointly optimizes them with personality prediction via multi-task learning. A cross-attention module is employed to capture fine-grained interactions between personality traits and the inferred emotional representations. To further refine relational reasoning, EmoPerso adopts a self-taught strategy to enhance the model's reasoning capabilities iteratively. Extensive experiments on two benchmark datasets demonstrate that EmoPerso surpasses state-of-the-art models. The source code is available at <https://github.com/slz0925/EmoPerso>.

## CCS Concepts

• **Computing methodologies** → **Information extraction; Multi-task learning; Supervised learning by classification; Unsupervised learning; Knowledge representation and reasoning.**

## Keywords

Personality Detection; Emotion Modelling; Multi-Task Learning; Reasoning Chains; Self-Supervised Learning

## ACM Reference Format:

Lingzhi Shen, Xiaohao Cai, Yunfei Long, Imran Razzak, Guanming Chen, and Shoaib Jameel. 2025. EmoPerso: Enhancing Personality Detection with

Self-Supervised Emotion-Aware Modelling. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management (CIKM '25)*, November 10–14, 2025, Seoul, Republic of Korea. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3746252.3761247>

## 1 Introduction

Imagine browsing social media and coming across a post describing an experience on “How to Quickly Improve Your Social Skills?”. This post resonated widely, with thousands of likes and comments [5]. However, its expression can actually be rewritten in different styles based on the author's personality type and emotional state. For instance, an extroverted version might emphasize active communication and group interactions, whereas an introverted version could focus more on building deep connections in small social settings [36]. This phenomenon raises a fundamental question: How do personality and emotion influence the way text is expressed? Furthermore, can we leverage personality and emotion features extracted from posts to predict a user's personality type, such as the Myers-Briggs Type Indicator (MBTI)<sup>1</sup> [4].

In cognitive science, the relationship between emotion and personality has been widely studied [38]. Personality traits reflect an individual's long-term behavioural patterns, whereas emotions are expressions of short-term mental states [43]. As illustrated in Figure 1, personality and emotion operate on different cognitive timescales and interact to influence human behaviour. This perspective is also aligned with the Cognitive-Affective Personality System (CAPS) theory [25], one of the most widely cited frameworks in personality psychology, which conceptualizes personality as a stable disposition that organizes how emotional reactions vary across situations. To systematically study emotional responses, Ekman's Basic Emotion Theory [9] provides a foundational taxonomy, positing that human emotions can be categorized into seven fundamental types: joy, anger, sadness, fear, disgust, surprise, and contempt. Building on this, studies suggest that individuals with different personality types may experience and express emotions



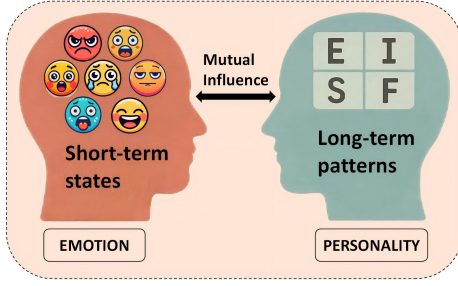
This work is licensed under a Creative Commons Attribution 4.0 International License. *CIKM '25, Seoul, Republic of Korea*

© 2025 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-2040-6/2025/11

<https://doi.org/10.1145/3746252.3761247>

<sup>1</sup>The MBTI is a widely used personality framework that classifies individuals into 16 types based on four dichotomies: Introversion (I) vs. Extraversion (E), Sensing (S) vs. Intuition (N), Thinking (T) vs. Feeling (F), and Perceiving (P) vs. Judging (J).



**Figure 1: An illustration of the cognitive distinction and interaction between emotional states and personality traits.**

differently when faced with the same situation [37, 52]. For example, extroverted individuals are more likely to exhibit joy and surprise, whereas introverted individuals may display a more reserved emotional response, such as subdued expressions of joy or a tendency toward introspective emotions like sadness or fear.

Recent approaches have attempted to incorporate emotional information into personality prediction to improve its accuracy. For example, Hu et al. [12] leverage large language models (LLMs) to generate augmented textual data and interpret personality labels from raw social media posts, focusing on semantic, sentimental, and linguistic aspects. While AI-generated data can provide valuable insights, these methods often lack targeted modelling of the relationship between personality and psycholinguistic factors, making it challenging for models to capture their fine-grained interactions [41]. Importantly, sentiment and emotion are distinct [14]: sentiment generally reflects coarse-grained polarity (e.g., positive or negative) without capturing the nuanced psychological states conveyed by emotions. Another model, EERPD [19], integrates emotion regulation with emotional features using few-shot learning and chain-of-thought reasoning. However, it may require high-quality and diverse emotional data, which can be difficult to obtain in real-world scenarios. Similarly, many existing methods heavily rely on large-scale labeled datasets [11, 40], yet acquiring high-quality personality-labeled data remains challenging [57], especially when input text is incomplete or noisy.

In this paper, we propose *EmoPerso*, a novel self-supervised **emotion-personality** joint learning framework, whose core idea is to infer pseudo-labeled emotion representations from personality-labeled social media posts to construct emotion-aware personality representations. To improve data diversity, *EmoPerso* incorporates LLM-based generative mechanisms [44], including style-conditioned paraphrasing and contextual feature completion. Inferred emotion signals serve as auxiliary supervision and are jointly optimized with personality prediction through multi-task learning (MTL) [13], allowing the model to capture low-level sharing between emotional and personality-related features. Furthermore, the framework introduces a cross-attention module to capture fine-grained personality modulation conditioned on emotion embeddings, reinforced by an emotion-conditioned weighting mechanism that enhances the representation of psychologically salient cues. In the reasoning stage, *EmoPerso* integrates the Self-Taught Reasoner (STaR) [53] to generate individualized reasoning chains, which are filtered using information gain and mutual information metrics

to further strengthen the semantic coupling between emotional features and personality traits.

**Key Contributions:** Inspired by cognitive science and Basic Emotion Theory, we propose *EmoPerso*, a novel self-supervised emotion-personality joint learning framework, which verifies the importance of emotional features for personality detection tasks. Unlike traditional machine learning studies that treat personality prediction and emotion analysis as independent tasks while ignoring their interplay, *EmoPerso* leverages LLM-driven generative mechanisms, multi-task learning, cross-attention modelling, and enhanced reasoning chains, through which the interaction between pseudo-labelled emotion signals and personality traits is progressively deepened, enabling the model to gradually construct emotion-aware personality predictions. This unified strategy refines the learning process by utilizing self-generated data, introducing a flexible text augmentation paradigm that reduces reliance on external annotations. Extensive experiments demonstrate that *EmoPerso* outperforms state-of-the-art models on benchmark datasets.

## 2 Related Work

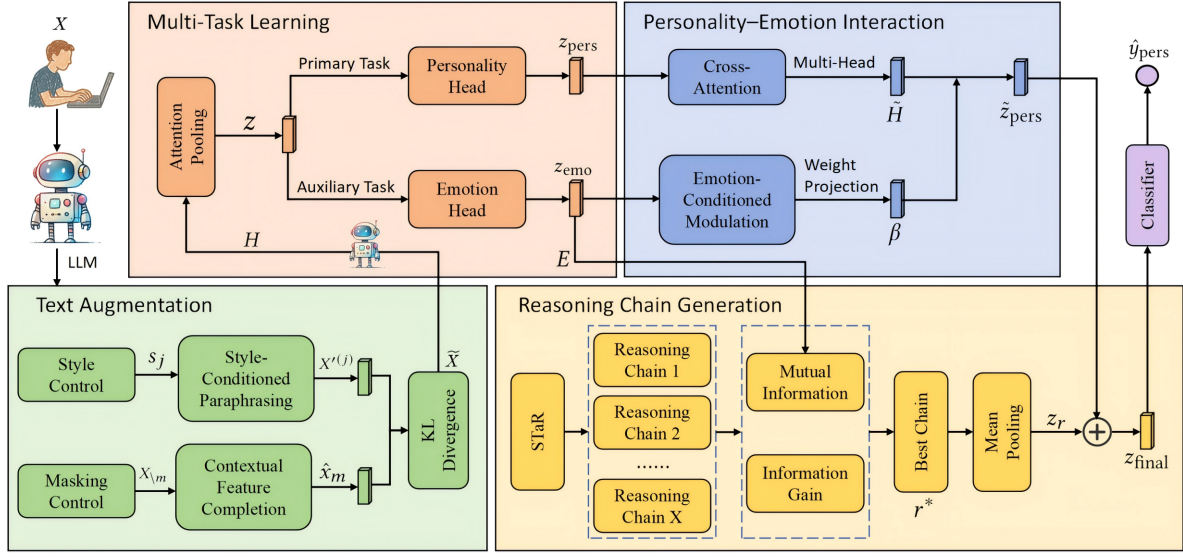
Previous research on personality detection spans a wide range of model architectures and learning paradigms. To provide a clear and comparable overview, we categorize existing methods into two main groups: classic deep learning methods and LLM-driven methods. This categorization not only reflects the methodological evolution of the field, but also highlights the shift in research focus toward cognitive reasoning.

### 2.1 Classical Deep Learning Methods

Early neural networks laid the foundation for modelling sequential and structural dependencies in user-generated content [34]. LSTM [29] and GRU [31] architectures are frequently employed to capture temporal and spatial dependencies in personality traits. Some studies have adopted hierarchical feature extraction approaches [6, 39], while Transformer-based models leverage self-attention mechanisms to capture long-range dependencies, thereby improving their ability to model global relationships [23]. Additionally, GNNs have been utilized to model the complex interactions between personality traits and external factors [33, 58]. However, most deep learning methods, trained in supervised settings, rely on fixed features or pretrained data [59], limiting adaptability and generalizability.

### 2.2 LLM-Driven Methods

Recent advances in LLMs have provided enhanced generalisation and reasoning capabilities, making them increasingly popular across a wide range of downstream tasks, including personality detection [50]. Studies have shown [3] that LLMs can accurately classify personality traits based on social media data and outperform traditional machine learning methods through prompt engineering and few-shot learning techniques [26]. Instruction-tuned models, such as GPT-4o [1] and Claude 3.5 Sonnet [27], have improved the reliability of personality assessment tools by generating self-evaluative texts aligned with psychological frameworks [30], while Llama [10] has demonstrated strong performance in open-domain settings [32, 56]. However, these methods still largely depend on shallow pattern matching rather than genuinely comprehending



**Figure 2: Overall architecture of EmoPerso. The framework leverages LLMs for self-supervised emotion feature extraction. It integrates LLM-based generative mechanisms, MTL, cross-attention, and STaR to enhance personality-emotion interactions.**

the logical structure of information in psychological theory reasoning tasks [17]. They also struggle to identify key information in subtle cognitive reasoning tasks [2].

Unlike prior studies that rely on prompt engineering or externally annotated emotion data, our EmoPerso introduces a self-supervised learning framework that infers emotion representations from personality texts using pseudo-labels. These emotion features are jointly optimized with personality traits in learning process, enabling the model to internalize emotion-personality dependencies. In addition, EmoPerso incorporates LLM-based data synthesis techniques to improve generalization and employs reasoning chains to enhance its inference quality.

### 3 Our Novel EmoPerso Framework

The design of EmoPerso (Figure 2 and Algorithm 1) introduces a novel self-supervised joint learning framework that leverages LLMs to improve personality detection through emotion-aware modelling. It enhances text generalization via LLM-based generation mechanisms, jointly optimizes emotion and personality representations through MTL, and refines their fine-grained interactions via cross-attention. Finally, STaR is employed to generate and select informative reasoning chains, further enhancing emotion-conditioned personality inference.

#### 3.1 Text Augmentation

We leverage large language models (LLMs) to perform self-supervised text augmentation that enhances representation diversity. Specifically, we introduce two complementary mechanisms: style-conditioned paraphrasing and contextual feature completion. Both are designed to generate semantically consistent yet informationally diverse variants of the input, facilitating better generalisation under limited labelled data.

#### Algorithm 1 Pseudocode of EmoPerso

**Require:** Dataset  $\mathcal{D}$ , training epochs  $N$ , batch size  $B$ , and learning rate  $\eta$

**Ensure:** Optimized parameters  $\Theta$

- 1: Initialize optimizer and loss weights  $\lambda_{\text{MTL}}, \lambda_{\text{cross}}, \lambda_{\text{star}}$
- 2: **for** epoch in range( $N$ ) **do**
- 3:   Load batch  $(X, y_{\text{pers}})$
- 4:    $\tilde{X} \leftarrow \text{AUGMENT}(X)$
- 5:    $z_{\text{shared}} \leftarrow \text{ENCODE}(\tilde{X})$
- 6:    $(z_{\text{pers}}, z_{\text{emo}}) \leftarrow \text{DECOMPOSE}(z_{\text{shared}})$
- 7:    $\tilde{z}_{\text{pers}} \leftarrow \text{INTERACT}(z_{\text{pers}}, z_{\text{emo}}, \tilde{X})$
- 8:    $R \leftarrow \text{GENERATECHAINS}(X)$
- 9:    $\{P(r_i)\} \leftarrow \text{SCORE}(R; z_{\text{emo}})$
- 10:    $r^* \leftarrow \text{SELECT}(R; \{P(r_i)\})$
- 11:    $z_r \leftarrow \text{EMBED}(r^*)$
- 12:    $z_{\text{final}} \leftarrow \text{COMBINE}(\tilde{z}_{\text{pers}}, z_r)$
- 13:    $\hat{y}_{\text{pers}} \leftarrow \text{INFER}(z_{\text{final}})$
- 14:    $\mathcal{L}_{\text{MTL}} \leftarrow \mathcal{L}_{\text{pers}} + \mathcal{L}_{\text{emo}}$
- 15:    $\mathcal{L}_{\text{total}} \leftarrow \lambda_{\text{MTL}} \mathcal{L}_{\text{MTL}} + \lambda_{\text{cross}} \mathcal{L}_{\text{cross}} + \lambda_{\text{star}} \mathcal{L}_{\text{star}}$
- 16:    $\Theta \leftarrow \text{UPDATE}(\Theta; \nabla \mathcal{L}_{\text{total}})$
- 17: **end for**

Given an input post  $X = (x_1, x_2, \dots, x_T)$ , where  $x_i \in \mathbb{V}$  and  $\mathbb{V}$  denotes the vocabulary space, we generate stylistically diverse paraphrases conditioned on a control signal  $s_j$  indicating attributes such as formality, expressiveness, or conciseness. These three styles are selected for their strong alignment with core personality dimensions, their prevalence in social media discourse, and their ability to guide generation in semantically meaningful and distinguishable ways. The LLM acts as a conditional generator, producing the  $j$ -th paraphrased variant as  $X^{(j)} = \text{LLM}(X \mid s_j)$ , where  $j = 1, \dots, k$



and  $k$  is the total number of style conditions. This process relies solely on the prompt-driven capabilities of the LLM backbone.

To encourage diversity among paraphrases, we adopt sampling-based decoding strategies, specifically nucleus sampling (i.e., top- $p$  sampling) [21], instead of deterministic generation [7], which (e.g., greedy search) tends to produce repetitive and generic outputs that lack stylistic variation, limiting the model’s ability to explore diverse surface realisations of personality-related content. This decoding strategy is natively supported by LLMs and is invoked during generation to promote lexical and stylistic variability.

To handle incomplete or noisy inputs, we simulate missing information by randomly masking spans in  $X$ , resulting in  $X_{\setminus m}$ , where  $m \subseteq \{1, \dots, T\}$  indicates masked positions. The LLM is then prompted to reconstruct the missing content based on surrounding context, predicting the masked tokens as  $\hat{x}_m = \text{LLM}(X_{\setminus m})$ . Similar to paraphrasing, contextual completion also leverages the inherent infilling capability of the LLM. In practice, the masked spans are sampled at the phrase or content-word level rather than the token level to ensure syntactic and semantic plausibility.

To ensure semantic fidelity between original and augmented sequences, we impose a regularisation loss based on Kullback-Leibler (KL) divergence [8] between their token-level output distributions. Let  $P(x_t | X)$  and  $P(x_t | \hat{X})$  denote the predicted token distributions at position  $t$  from the LLM, conditioned on the original and augmented inputs, respectively. KL divergence is computed over aligned token positions as:  $\mathcal{L}_{\text{KL}} = \frac{1}{T} \sum_{t=1}^T D_{\text{KL}}(P(x_t | X) \| P(x_t | \hat{X}))$ . This regularisation assumes access to the LLM’s token-level softmax probabilities, feasible in open-source implementations.

In practice, we access token distributions via the LLM’s softmax outputs at each generation step. To account for stylistic variations, we also include a supervised style classification loss  $\mathcal{L}_{\text{style}}$  using pseudo-labels corresponding to each style condition. This component is implemented using a lightweight two-layer MLP trained jointly with the main objectives. The overall augmentation loss is thus defined as:  $\mathcal{L}_{\text{gen}} = \lambda_{\text{style}} \mathcal{L}_{\text{style}} + \lambda_{\text{KL}} \mathcal{L}_{\text{KL}}$ , where  $\lambda_{\text{style}}$  and  $\lambda_{\text{KL}}$  are balancing coefficients.

This augmentation strategy not only mitigates data sparsity but also exposes the model to emotionally and stylistically conditioned variants, which are critical for capturing fine-grained personality signals. Example prompts and style control templates are provided in Figure 3. Given a fixed input sentence and a list of target styles, the model dynamically constructs prompts and generates stylistically diverse outputs using an LLM.

### 3.2 Multi-Task Learning

We adopt an LLM as a self-supervised feature extractor to jointly capture signals related to both personality and emotion, enabling MTL over shared latent representations. Unlike conventional approaches that treat personality and emotion as separate tasks trained on annotated datasets [18], the core self-supervised property of our framework lies in the fact that emotion labels are not externally provided. Instead, emotion signals are inferred directly from personality-labeled data through auxiliary modelling, without any explicit supervision. The emotion head is trained using pseudo-labels derived from stylistic and emotional cues in the input, such

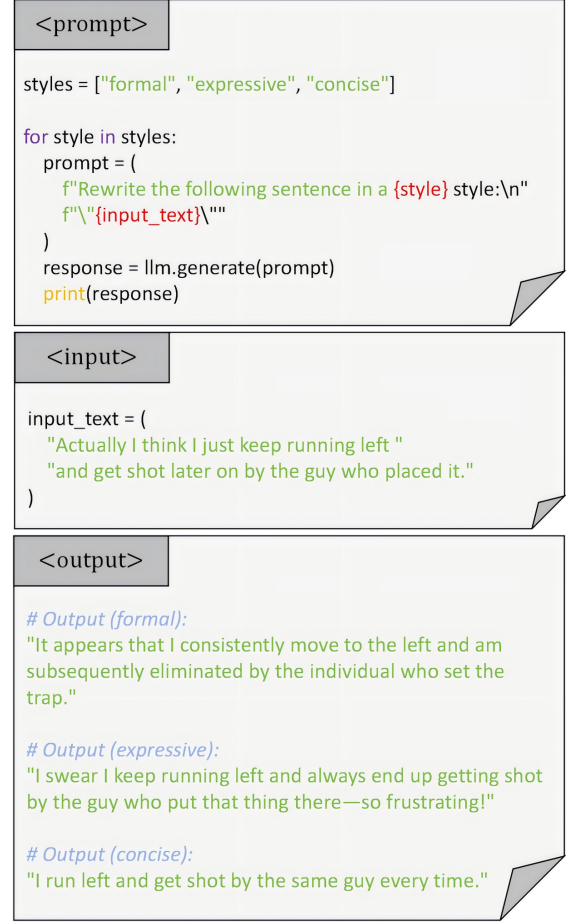


Figure 3: Code-style illustration of style-conditioned paraphrasing using a real example from the Kaggle dataset.

as lexical choice, punctuation usage, and emotional valence intensity. These automatically generated signals serve as supervisory targets, allowing EmoPerso to construct emotion-aware representations in a purely self-supervised manner throughout the entire training process.

Given an augmented input post  $\tilde{X} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_T)$ , produced via style transformation and feature completion, the LLM encodes the sequence into a hidden representation  $H = (h_1, h_2, \dots, h_T)$ , where  $h_i \in \mathbb{R}^d$  denotes the embedding of token  $\tilde{x}_i$ .

To aggregate token-level information into a global representation, we apply an attention pooling mechanism over  $H$ , i.e.,

$$z = \sum_{i=1}^T \alpha_i h_i, \quad \alpha_i = \frac{\exp(v^\top \tanh(W h_i + b))}{\sum_{j=1}^T \exp(v^\top \tanh(W h_j + b))}, \quad (1)$$

where  $W \in \mathbb{R}^{d' \times d}$ ,  $b \in \mathbb{R}^{d'}$ , and  $v \in \mathbb{R}^{d'}$  are trainable parameters. Technically, this attention pooling mechanism is implemented as a single-head feedforward scoring function using a two-layer MLP followed by softmax normalisation. This design enhances the model’s ability to focus on semantically salient tokens, which is

especially important for personality and emotion modelling, where relevant signals are often sparse and context-dependent.

We treat personality detection as four independent binary classification tasks, corresponding to the MBTI dimensions: Introversion/Extraversion (I/E), Sensing/Intuition (S/N), Thinking/Feeling (T/F), and Perceiving/Judging (P/J). Let  $C_p = 4$  be the number of dimensions. Both personality and emotion tasks project the shared representation  $z$  into task-specific embeddings using lightweight two-layer MLPs, denoted as  $z_{\text{pers}} = f_{\text{pers}}(z)$  and  $z_{\text{emo}} = f_{\text{emo}}(z)$ , where  $z_{\text{pers}}, z_{\text{emo}} \in \mathbb{R}^d$  represent personality-specific and emotion-specific representations, respectively. These embeddings are used to produce the corresponding classification logits for each task.

The personality prediction head outputs a 4-dimensional logit vector, and the probability for each dimension is computed using sigmoid activation as  $p_c = \sigma((W_{\text{pers}} z_{\text{pers}} + b_{\text{pers}})_c)$ , where  $c = 1, \dots, C_p$ ,  $W_{\text{pers}} \in \mathbb{R}^{C_p \times d}$ , and  $b_{\text{pers}} \in \mathbb{R}^{C_p}$  are learnable parameters. The corresponding binary cross-entropy loss is:

$$\mathcal{L}_{\text{pers}} = - \sum_{c=1}^{C_p} (y_c \log p_c + (1 - y_c) \log(1 - p_c)), \quad (2)$$

where  $y_c \in \{0, 1\}$  is the ground-truth label for the  $c$ -th personality dimension.

Emotion prediction is modelled as a multi-label classification task over  $C_e$  emotion categories. Importantly, these emotion labels are not provided by external annotation. Instead, the model learns to predict emotion categories based on latent emotional cues that co-occur in personality-labeled text. These cues include stylistic and psycholinguistic markers such as valence-bearing adjectives (e.g., “excited,” “frustrated”), intensifiers (e.g., “really,” “extremely”), exclamation usage, emotive punctuation, and affective n-grams, which are often indicative of underlying emotional states. This enables the emotion stream to act as a self-supervised auxiliary task. The prediction is computed using sigmoid activation as  $\tilde{p}_c = \sigma((W_{\text{emo}} z_{\text{emo}} + b_{\text{emo}})_c)$ , where  $c = 1, \dots, C_e$ ,  $W_{\text{emo}} \in \mathbb{R}^{C_e \times d}$ , and  $b_{\text{emo}} \in \mathbb{R}^{C_e}$  are trainable parameters. The corresponding multi-label binary cross-entropy loss is:

$$\mathcal{L}_{\text{emo}} = - \sum_{c=1}^{C_e} (y_c \log \tilde{p}_c + (1 - y_c) \log(1 - \tilde{p}_c)), \quad (3)$$

where  $y_c \in \{0, 1\}$  denotes the pseudo-label for the  $c$ -th emotion category, selected from the inferred label set  $\hat{Y}_{\text{emo}}$  automatically constructed based on emotional cues in the input.

The overall MTL objective jointly optimises both tasks, i.e.,

$$\mathcal{L}_{\text{MTL}} = \lambda_{\text{pers}} \mathcal{L}_{\text{pers}} + \lambda_{\text{emo}} \mathcal{L}_{\text{emo}}, \quad (4)$$

where  $\lambda_{\text{pers}}$  and  $\lambda_{\text{emo}}$  are hyperparameters balancing the two tasks.

### 3.3 Personality–Emotion Interaction

Personality traits and emotional states, while related, often exhibit asymmetric and context-dependent correlations at the token level [28]. For example, emotional expressions tend to occur sparsely in text but have varying influences across different personality dimensions. To realize this intuition, we design a two-stage interaction mechanism: a multi-head cross-attention module that aligns personality with token-level input, and an emotion-conditioned modulation layer that reweights token contributions based on emotional

context. This design enables the model to align emotionally salient input regions with personality-relevant features, going beyond standard attention layers that treat auxiliary signals as uniformly distributed context [16].

Specifically, to enhance fine-grained interaction modelling between personality and emotion traits, we refine the personality-specific representation  $z_{\text{pers}}$ , which is obtained from the MTL head, via a multi-head cross-attention mechanism. This step enables the model to selectively attend to relevant contextual tokens based on personality semantics, while later incorporating emotion-conditioned modulation to further refine the final representation.

Given the hidden token sequence  $H = (h_1, h_2, \dots, h_T) \in \mathbb{R}^{T \times d}$  and the personality-specific query vector  $z_{\text{pers}} \in \mathbb{R}^d$ , we compute multi-head cross-attention by projecting  $z_{\text{pers}}$  as a query and  $H$  as keys and values. Specifically, for each head  $h \in \{1, \dots, H\}$ , the query, key, and value matrices are computed as:

$$Q^{(h)} = W_Q^{(h)} z_{\text{pers}}, \quad K^{(h)} = W_K^{(h)} H, \quad V^{(h)} = W_V^{(h)} H, \quad (5)$$

where  $Q^{(h)} \in \mathbb{R}^{1 \times d_k}$ ,  $K^{(h)}, V^{(h)} \in \mathbb{R}^{T \times d_k}$ , and all  $W_*^{(h)}$  are head-specific trainable projection matrices. This attention mechanism is implemented using standard scaled dot-product attention, where the outputs are computed as:

$$A^{(h)} = \text{Softmax} \left( \frac{Q^{(h)} K^{(h)^\top}}{\sqrt{d_k}} \right) V^{(h)}. \quad (6)$$

The resulting attended token representations  $\tilde{H} = (\tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_T)$  are obtained by concatenating outputs from all heads and projecting them via a learned output matrix  $W_O \in \mathbb{R}^{H d_k \times d}$ , i.e., each token representation is computed as  $\tilde{h}_i = \text{Concat}(A_i^{(1)}, \dots, A_i^{(H)}) W_O$ .

To further incorporate emotion-specific context, we introduce an emotion-conditioned attention modulation. Based on the emotion embedding  $z_{\text{emo}} \in \mathbb{R}^d$ , predicted from pseudo-label-guided emotion supervision in the MTL stage, we compute token-level importance weights as  $\beta = \text{Softmax}(W_{\text{emo}} z_{\text{emo}})$ , where  $W_{\text{emo}} \in \mathbb{R}^{T \times d}$  is a learned projection matrix. These weights are used to aggregate the attended token representations such that the final personality vector is given by  $\tilde{z}_{\text{pers}} = \sum_{i=1}^T \beta_i \tilde{h}_i$ . This design allows the model to selectively amplify psychologically salient features from the diverse semantic subspaces constructed in the previous cross-attention step, guided by emotional context.

To align personality- and emotion-guided features, we introduce a consistency regularisation loss based on cosine similarity, defined as  $\mathcal{L}_{\text{cross}} = 1 - \cos(\tilde{z}_{\text{pers}}, z_{\text{emo}})$ , where  $\cos(\cdot, \cdot)$  denotes the cosine similarity between the final personality representation and the emotion embedding. This consistency loss does not simply align representations geometrically, but instead guides the model to preserve emotionally discriminative features within the personality space by minimizing the distance between the emotion embedding and the personality representation.

The final training objective integrates this regularisation with the multi-task loss:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{MTL}} \mathcal{L}_{\text{MTL}} + \lambda_{\text{cross}} \mathcal{L}_{\text{cross}}, \quad (7)$$

where  $\lambda_{\text{MTL}}$  and  $\lambda_{\text{cross}}$  are trade-off hyperparameters.

### 3.4 Reasoning Chain Generation

To further enhance the model’s reasoning capabilities, we adopt a self-taught rationale generation strategy inspired by STaR, encouraging the model to generate its intermediate reasoning chains. We extend this process by introducing an information-theoretic selection mechanism to identify the most informative reasoning chains, particularly those capturing relational cues between personality and emotion features.

For each input post  $X$ , we generate a set of candidate reasoning chains  $R = \{r_1, r_2, \dots, r_n\}$ , where each  $r_i = (s_1^{(i)}, s_2^{(i)}, \dots, s_L^{(i)})$  consists of a sequence of intermediate reasoning steps, with  $s_j^{(i)}$  denoting the  $j$ -th step and  $L$  the chain length. We use an LLM decoder to autoregressively generate each reasoning step, employing nucleus sampling with temperature control to encourage diverse and coherent chains, i.e.,

$$s_j^{(i)} = \arg \max_s P(s | s_1^{(i)}, \dots, s_{j-1}^{(i)}, X), \quad (8)$$

where  $P(\cdot)$  denotes the token-level generation distribution.

Since not all reasoning chains contribute equally to personality prediction, we apply a filtering mechanism based on information-theoretic criteria. The information gain (IG) [35] measures how much uncertainty is reduced in personality classification given a reasoning chain  $IG(r_i) = H(Y) - H(Y | r_i)$ , where  $H(\cdot)$  denotes the entropy of the predicted label distribution. In practice, both terms are approximated using the model’s predicted log-probabilities with and without conditioning on the reasoning chain  $r_i$ .

Additionally, to capture the strength of emotion–personality coupling, we compute the mutual information (MI) [45] between each reasoning chain and the extracted emotion features, i.e.,

$$MI(r_i, E) = \sum_{y \in Y} \sum_{e \in E} P(y, e) \log \frac{P(y, e)}{P(y)P(e)}, \quad (9)$$

where  $E$  denotes the emotion feature set extracted via self-supervised learning, conditioned on pseudo-labels from the MTL module. Joint and marginal distributions  $P(y, e)$  are estimated over mini-batches using empirical frequency counts from predicted outputs. We define a normalised preference score over candidate chains as:

$$P(r_i) = \frac{\exp(\lambda_{IG} IG(r_i) + \lambda_{MI} MI(r_i, E))}{\sum_{j=1}^n \exp(\lambda_{IG} IG(r_j) + \lambda_{MI} MI(r_j, E))}. \quad (10)$$

These preference scores are dynamically updated during training and used both for chain selection and for regularisation. The optimal reasoning chain is selected by maximising the combined signal:

$$r^* = \arg \max_{r_i \in R} (\lambda_{IG} IG(r_i) + \lambda_{MI} MI(r_i, E)), \quad (11)$$

where  $\lambda_{IG}$  and  $\lambda_{MI}$  are balancing coefficients.

To incorporate the selected reasoning chain into the model, we first encode  $r^*$  into a vector representation. Specifically, we tokenize  $r^*$ , pass it through the LLM, and apply mean pooling over its token embeddings to obtain a reasoning embedding  $z_r \in \mathbb{R}^d$ . This vector captures the semantic content of the selected rationale in the same latent space as the emotion-aware personality representation  $\tilde{z}_{pers}$ .

We integrate the reasoning vector  $z_r$  with the personality-specific representation  $\tilde{z}_{pers}$  by applying a lightweight transformation to their concatenation, i.e.,  $z_{final} = \text{MLP}(\text{Concat}(\tilde{z}_{pers}, z_r))$ , which projects the fused representation back into  $\mathbb{R}^d$ . This design allows

**Table 1: Statistics on the quantity and class distribution for the Kaggle and Pandora datasets.**

Dataset	Types	Train	Validation	Test
Kaggle	I/E	1194 / 4011	409 / 1326	396 / 1339
	S/N	610 / 4478	222 / 1513	248 / 1487
	T/F	2410 / 2795	791 / 944	780 / 955
	P/J	2109 / 3096	672 / 1063	653 / 1082
	Posts	246794	82642	82152
Pandora	I/E	1162 / 4278	386 / 1427	377 / 1437
	S/N	727 / 4830	208 / 1605	210 / 1604
	T/F	3549 / 1891	1120 / 693	1182 / 632
	P/J	2229 / 3211	770 / 1043	758 / 1056
	Posts	523534	173005	174080

the model to incorporate high-level inductive signals from the reasoning chain into the personality representation, while preserving the original semantic structure derived from multi-task learning and attention-based interactions.

To encourage confident reasoning selection, we define a reasoning chain entropy loss  $\mathcal{L}_{star} = -\sum_{i=1}^n P(r_i) \log P(r_i)$ , which penalizes overly uniform distributions over candidate rationales.

The final training objective combines all loss components, i.e.,

$$\mathcal{L}_{total} = \lambda_{MTL} \mathcal{L}_{MTL} + \lambda_{cross} \mathcal{L}_{cross} + \lambda_{star} \mathcal{L}_{star}, \quad (12)$$

where each  $\lambda(\cdot)$  controls the contribution of its corresponding loss term. The resulting fused representation  $z_{final}$  is then used to compute the final personality prediction  $\hat{y}_{pers}$ .

## 4 Experiments and Results

This section details the benchmark datasets used, experimental settings, and the significant results obtained to determine whether EmoPerso outperforms recent robust models through rigorous evaluation of the effectiveness of EmoPerso on personality detection tasks. We further perform a comprehensive ablation study and visualisation-based analysis to assess the contribution of each component, followed by an in-depth qualitative analysis.

### 4.1 Experimental Setup

**Datasets:** To ensure a fair comparison with previous work, we selected the same two datasets, i.e., Kaggle<sup>2</sup> and Pandora<sup>3</sup>. The Kaggle dataset is sourced from the PersonalityCafe forum, an online community focused on discussions about personality types. This dataset contains posts from 8,675 users, with each user contributing approximately 45 to 50 posts. The posts cover a variety of topics, including psychology, personal experiences, and everyday discussions. The dataset is labelled according to users’ self-reported MBTI personality dimensions. Pandora is a larger corpus from the Reddit platform, which includes MBTI labels for 9,084 users, primarily extracted from the flairs (short self-descriptions) in MBTI-related subreddits. The number of posts per user ranges from dozens to hundreds, and due to the diversity of the Reddit community, the content covers a broader range of topics. Table 1 presents some statistics of the two datasets.

**Implementation Details:** We use a frozen DeepSeek-V3 [20] as the backbone. To ensure cross-task consistency despite differing

<sup>2</sup><https://www.kaggle.com/datasnaek/mbti-type>

<sup>3</sup><https://psy.takelab.fer.hr/datasets/all>

**Table 2: Comparison of EmoPerso with state-of-the-art baselines on the Kaggle and Pandora datasets. Results are reported using Macro-F1 (%) scores across the four MBTI dimensions and the overall average (Avg).**

Methods	Kaggle					Pandora				
	I/E	S/N	T/F	P/J	Avg	I/E	S/N	T/F	P/J	Avg
AttRCNN	59.74	64.08	78.77	66.44	67.25	48.55	56.19	64.39	57.26	56.60
SN+Attn	65.43	62.15	78.05	63.92	67.39	56.98	54.78	60.95	54.81	56.88
Transformer-MD	66.08	69.10	79.19	67.50	70.47	55.26	58.77	69.26	60.90	61.05
PQ-Net	68.94	67.65	79.12	69.57	71.32	57.07	55.26	65.64	58.74	59.18
TrigNet	69.54	67.17	79.06	67.69	70.86	56.69	55.57	66.38	57.27	58.98
PS-GCN	70.52	65.73	70.51	67.13	68.47	59.12	54.88	67.35	58.62	59.49
D-DGCN	69.52	67.19	80.53	68.16	71.35	59.98	55.52	70.53	59.56	61.40
DEN	69.95	66.39	80.65	69.02	71.50	60.86	57.74	71.64	59.17	62.35
MvP	67.68	69.89	80.99	68.32	71.72	60.08	56.99	69.12	61.19	61.85
PsyCoT	66.56	61.70	74.80	57.83	65.22	60.91	57.12	66.45	53.34	59.45
TAE	70.90	66.21	81.17	70.20	72.07	62.57	61.01	69.28	59.34	63.05
<b>EmoPerso</b>	<b>80.05</b>	<b>79.27</b>	<b>87.03</b>	<b>77.91</b>	<b>81.07</b>	<b>66.84</b>	<b>68.15</b>	<b>71.90</b>	<b>67.51</b>	<b>68.60</b>

token lengths, the input sequence is standardized to 2,048 tokens (median) with a hidden size of 4,096. Training employs Adam with learning rate  $1 \times 10^{-3}$ . For augmentation, we generate  $k = 3$  diverse paraphrases per input using prompt-based top- $p$  sampling ( $p = 0.9$ , temperature 1.0), capped at 512 tokens. Contextual completion masks 10% of tokens, generating up to 20 per span. KL divergence regularization is weighted by 0.1. For MTL, the loss ratio between personality and emotion tasks is 0.7:0.3, optimized with binary cross-entropy, reflecting the primary role of personality prediction and the auxiliary role of emotion modelling. Pseudo-labels for emotion are derived from affective heuristics (e.g., adjectives, intensifiers, and punctuation). Classifier heads are two-layer MLPs with ReLU and dropout 0.2. The cross-attention module uses four heads with residual connections and layer normalization. Reasoning chains are generated by the LLM ( $\leq 4$  steps), scored by information gain and mutual information from chain-conditioned vs. unconditioned probabilities. Preference scores are computed dynamically with softmax for selection and consistency loss. To avoid leakage, words or phrases directly matching personality labels are removed in preprocessing. Data is split 60/20/20 for train/validation/test, with results averaged over ten runs. Training is conducted on a cluster of NVIDIA H200 GPUs.

**Evaluation Metrics:** Macro-F1 has been widely used in previous studies and has become the standard evaluation metric for this task. Therefore, we adhere to this convention and use Macro-F1 to ensure consistency with prior work [12, 49, 55].

**Comparative Models:** We selected diverse sophisticated architectures. AttRCNN [46] integrates attention into an RCNN structure with a CNN-Inception module for robust feature extraction. SN+Attn [24] uses a Sequence Network with dual attention at message and word levels to enhance signal relevance. Transformer-MD [47] employs Transformer-XL and memory mechanisms for disorder-agnostic post integration with dimension-specific attention. PQ-Net [48] fuses psychological questionnaires and user text via cross-attention to capture personality cues. TrigNet [51] constructs a heterogeneous tripartite graph with flow graph attention (GAT) for psycholinguistic integration. PS-GCN [22] merges psycholinguistic knowledge graphs and sentiment semantics via GCN and multi-head attention. D-DGCN [49] dynamically builds graph structures, integrating multiple posts disorder-agnostically. DEN

[54] models long-term personality traits with GCN, short-term states with BERT, and enhances both via bidirectional interaction. MvP [55] employs a multi-view Mixture-of-Experts with consistency regularization to integrate diverse perspectives of user posts. PsyCoT [50] structures psychological questionnaires as a reasoning chain, using multi-turn dialogue prompting for LLM-based scoring. TAE [12] leverages LLM-generated text augmentation and label explanations, applying contrastive learning to improve psycholinguistic representation.

## 4.2 Overall Results

The comparison of Macro-F1 scores between our EmoPerso and the baseline models is presented in Table 2. EmoPerso achieved the highest performance across all four dimensions as well as the overall average. On the Kaggle dataset, EmoPerso attained an average score of 81.07%, outperforming the best existing model, TAE, by 9.00%. This marks the first time that a personality detection model has surpassed 80% on this dataset, establishing a new milestone for future research. Similarly, both the I/E and S/N dimensions have reached approximately 80% for the first time. On the Pandora dataset, EmoPerso achieved an average score of 68.60%, surpassing TAE by 5.55%. Notably, compared to the baselines, EmoPerso demonstrated significant improvements under severe class imbalance, particularly reducing the gap between T/F and the other three dimensions.

## 4.3 Ablation Study

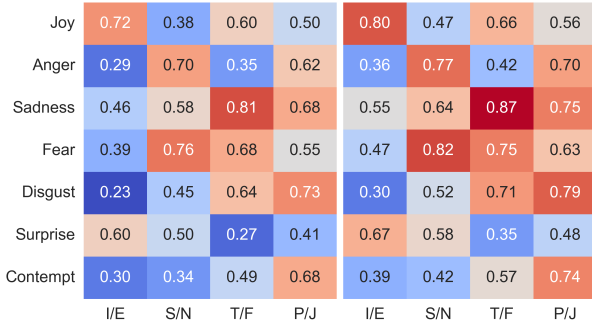
The results of the ablation study are shown in Table 3. First, we evaluate the performance of Vanilla DeepSeek-V3, which does not incorporate any additional optimizations. The results show a significant drop across all four dimensions and the overall average on both the Kaggle and Pandora datasets. This suggests that the designs introduced in EmoPerso are both necessary and effective, systematically addressing the base model’s limitations in emotional understanding and interaction-driven reasoning.

Next, we remove the emotion features, a key inspiration behind EmoPerso’s design. This results in a substantial drop of 8.04% on Kaggle and 5.39% on Pandora, making it the most impactful component in our ablation study. The reason for this is that emotion information serves as a crucial auxiliary signal for EmoPerso, and

**Table 3: Ablation results of different component configurations in EmoPerso. Reporting the Macro-F1 scores (%) on Kaggle and Pandora datasets for four dimensions and overall average (Avg).**

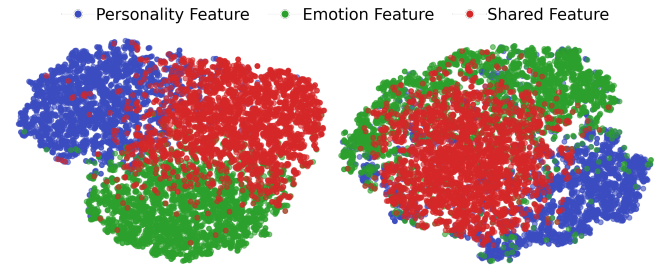
Components	Kaggle					Pandora				
	I/E	S/N	T/F	P/J	Avg	I/E	S/N	T/F	P/J	Avg
Vanilla DeepSeek-V3	65.16	59.75	77.71	64.43	66.76	60.17	55.93	65.58	60.35	60.51
w/o Emotions	69.31	73.79	76.01	73.00	73.03	58.85	65.69	63.15	65.16	63.21
w/o Generative Mechanism	74.43	71.22	84.30	76.11	76.52	61.61	60.90	70.16	66.41	64.77
w/o Paraphrasing	77.21	74.03	85.11	76.88	78.31	64.08	65.01	70.92	66.74	66.69
w/o KL Divergence	79.41	78.12	86.41	77.03	80.24	66.02	66.97	71.33	66.45	67.69
w/o MTL	75.60	71.21	78.08	74.48	74.84	63.44	61.91	64.78	65.69	63.96
w/o Shared Encoder	76.92	76.83	83.55	75.01	78.58	64.03	65.81	69.01	65.16	66.50
w/o CrossAttn Mechanism	76.47	77.05	78.94	71.95	76.10	63.96	66.83	65.14	62.10	64.51
Replace CrossAttn with Gated Fusion	78.32	77.93	82.36	74.88	78.87	65.15	67.32	68.03	63.97	66.12
w/o Emotion Modulation	76.83	77.14	80.91	74.48	77.84	64.12	66.58	68.21	63.88	65.70
w/o Reasoning Chains	74.70	75.03	78.82	73.15	75.42	63.04	64.78	65.51	63.27	64.15
Replace STaR with CoT Templates	76.43	76.81	81.90	75.56	77.68	64.70	65.85	68.10	64.98	65.91
w/o IG and MI	78.73	77.12	85.41	75.58	79.21	65.47	66.15	70.69	66.22	67.13
Replace DeepSeek-V3 with GPT-4o	<b>81.12</b>	78.23	84.45	<b>79.21</b>	80.75	<b>68.45</b>	67.98	70.10	<b>69.77</b>	<b>69.08</b>
<b>EmoPerso</b>	80.05	<b>79.27</b>	<b>87.03</b>	77.91	<b>81.07</b>	66.84	<b>68.15</b>	<b>71.90</b>	67.51	68.60

its absence deprives the model of its most essential enhancement. As shown in Figure 4, we use a heatmap to quantify the impact of different emotions on MBTI dimensions. The importance score for each emotion is derived from two sources: (1) its influence on the emotion-conditioned attention weights during personality representation refinement, and (2) its contribution to the selection of reasoning chains, measured by the integrated information gain and mutual information scores. The heatmap visualizes how different emotions differentially affect the four MBTI dimensions.

**Figure 4: Emotion contribution to prediction on the Kaggle (left) and Pandora (right) datasets.**

Removing the entire generative mechanism, including style-conditioned paraphrasing and contextual feature completion, leads to a substantial performance drop of 4.55% on the Kaggle dataset and 3.83% on the Pandora dataset. This result confirms the importance of personality-related text augmentation in addressing data sparsity and enhancing semantic diversity. When the style-conditioned paraphrasing component alone is removed, the performance decreases by 2.76% on Kaggle and 1.91% on Pandora. This finding highlights the contribution of paraphrasing to constructing more adaptive and expressive individual communication styles. The impact of removing KL divergence regularization is minimal. In its absence, slight semantic drift [15] is observed due to the lack of alignment between original and augmented sequences.

Following this, we eliminated the MTL designed to optimise emotion-personality interaction, which resulted in the loss of shared representations and the benefits of joint optimisation. Figure 5 uses t-SNE on the output representations from the personality head, emotion head, and shared latent vector after multi-task training, revealing substantial overlap between emotion and personality features and thus indicating rich shared representations. The visualization shows that MTL effectively captures and reinforces shared patterns, resulting in more cohesive clustering and better discriminability in the latent space. Additionally, we replaced the shared encoder with two separate encoders for the personality and emotion tasks. This modification introduces a hard separation between the two representation spaces, causing the model to lose the ability to transfer low-level linguistic and emotional cues across tasks.

**Figure 5: The t-SNE projection visualizes the shared features between personality and emotion under MTL, tested on the Kaggle (left) and Pandora (right).**

To evaluate the effect of fine-grained interaction modelling between personality and emotion, we first removed the cross-attention mechanism. This results in a significant performance drop of 4.97% on Kaggle and 4.09% on Pandora, confirming that enabling the personality representation to re-attend over the token sequence is critical for integrating emotion-conditioned cues. We then replaced the cross-attention mechanism with a gated fusion strategy, where token representations are modulated using a global gating vector



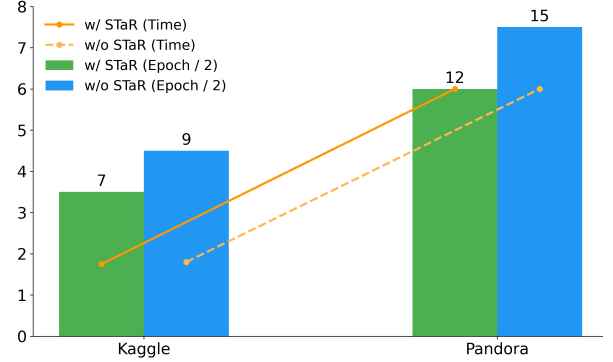
derived from both personality and emotion signals. While this alternative may reduce computational complexity, it leads to moderate performance decline. The result suggests that although gated fusion partially retains task-level interaction, it lacks the token-level selectivity and subspace diversity provided by multi-head attention. Finally, we removed the emotion modulation component leads to a performance drop of 3.23% on Kaggle and 2.90% on Pandora. These findings highlight the role of emotion-guided token weighting in dynamically amplifying psychologically salient cues.

In the case of reasoning chains, eliminating the STaR module also causes a noticeable performance drop, confirming its critical role in enhancing the model’s reasoning ability. Figure 6 compares the best training epoch (scaled) and total training time (in hours) for EmoPerso with and without STaR on the Kaggle and Pandora datasets. With STaR, the model achieves optimal performance in fewer epochs while maintaining comparable total training time. This suggests that STaR not only enhances personality inference through deeper reasoning but also accelerates the convergence process by guiding the model toward more informative and abstract patterns. To further probe the quality of reasoning, we replaced STaR with manually crafted Chain-of-Thought (CoT) templates [42], which are fixed prompting patterns designed to elicit step-by-step personality-related reasoning. For example, given a post, a CoT template might produce a generic rationale such as: “The user expresses frustration, which suggests emotional sensitivity, and therefore may lean toward the Feeling (F) trait.” This substitution results in moderate performance degradation compared to full STaR, but still performs better than completely removing the reasoning module. The result suggests that while CoT templates can provide basic interpretability, they fail to capture individualized, context-specific reasoning paths. Furthermore, when removing the IG and MI-based reasoning chain selection mechanism, the model still benefits from the existence of reasoning chains, but exhibits a moderate performance drop. This highlights that not all reasoning chains contribute equally to personality inference, and that selecting chains carrying the most informative relational signals is essential for fully leveraging the reasoning process.

Finally, we compare replacing DeepSeek-V3 with GPT-4o<sup>4</sup> as the backbone model and find their performance comparable. On Kaggle, DeepSeek-V3-based EmoPerso slightly outperforms GPT-4o-based EmoPerso on average, whereas on Pandora, the GPT-4o variant surpasses DeepSeek-V3 version. Interestingly, DeepSeek-V3 performs better on S/N and T/F, while GPT-4o excels on I/E and P/J, possibly because DeepSeek-V3 captures structured reasoning patterns, whereas GPT-4o better models conversational and social traits, directly influencing these dimensions.

## 5 Conclusion

This paper proposes EmoPerso, a novel self-supervised framework for joint emotion-personality modelling leveraging LLMs. By integrating generative mechanisms, MTL, and cross-attention, the framework facilitates deep interactions between personality and emotion, while STaR enhances the model’s emotion-conditioned



**Figure 6: Comparison of the best training epoch (scaled) and total training time with and without the STaR reasoning module across two datasets. Each bar indicates the optimal number of epochs (halved for visualization), and the lines represent total training time in hours.**

reasoning capabilities. Experiments on Kaggle and Pandora show EmoPerso surpasses state-of-the-art models, and ablation studies confirm the importance of each core component. EmoPerso also holds potential for real-world applications, such as mental health screening, bias-aware assessment, personalized marketing, and AI interaction systems. Future work includes multilingual adaptation, advanced generative reasoning, and modality extension.

## 6 Limitations and Ethical Considerations

While EmoPerso demonstrated promising outcomes, our experiments revealed several constraints. The current datasets suffer from label imbalance due to the scarcity of certain personality dimensions and biases in the data collection process. To enhance model’s generalization ability, more diverse and balanced real-world datasets are required, beyond the use of synthetically generated data alone. Moreover, the prevalence of social media bots further diminishes the reliability of the data.

Automated personality recognition entails a range of ethical challenges, including cultural bias, discrimination, and breaches of privacy. Strict adherence to ethical guidelines is essential to ensure fairness and safeguard data confidentiality. Most existing models are predominantly trained on English data, which significantly limits their applicability across different languages and cultural contexts. In addition, personality detection and classification based on emotions and spoken language are inherently complex and prone to biases against specific cultural or ethnic groups. The reliance on potentially biased data sources, such as self-reports and social media content, may undermine the credibility and practical value of applying the model in real-world scenarios.

## Acknowledgments

This work was supported by the Alan Turing Institute and Singapore’s DSO National Laboratories under a grant on improving multimodal misinformation detection through affective analysis.

<sup>4</sup>Note that although the closed nature of GPT-4o restricts components such as token-level KL divergence and self-supervised optimisation, these are approximated with inference-only strategies or replaced by compatible alternatives, enabling evaluation under a similar inference architecture.

## GenAI Usage Disclosure

This work employs DeepSeek-V3 as the backbone to perform controlled data augmentation (see Section 3.1), such as paraphrase generation and contextual feature completion, as part of our model's self-supervised training pipeline. These generative mechanisms were strictly used to enrich the training data and were evaluated for semantic consistency and task relevance. All model design, implementation, analysis, and writing decisions were made by the authors. No generative AI tools were listed as co-authors or assumed intellectual responsibility for any part of this work.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [2] Jaan Aru, Aqeel Labash, Oriol Corcoll, and Raul Vicente. 2023. Mind the gap: Challenges of deep learning approaches to theory of mind. *Artificial Intelligence Review* 56, 9 (2023), 9141–9156.
- [3] Goran Bubaš. 2024. The use of GPT-4o and Other Large Language Models for the Improvement and Design of Self-Assessment Scales for Measurement of Interpersonal Communication Skills. *arXiv preprint arXiv:2409.14050* (2024).
- [4] Sabur Butt, Grigori Sidorov, and Alexander Gelbukh. 2025. Interpretation of Myers–Briggs Type Indicator personality profiles based on ambivert continuum scale. *Expert Systems with Applications* 264 (2025), 125689.
- [5] Giovanni Luca Cascio Rizzo, Jonah Berger, Matteo De Angelis, and Rumen Pozharliev. 2023. How sensory language shapes influencer's impact. *Journal of Consumer Research* 50, 4 (2023), 810–825.
- [6] Quan Cheng and Wenwan Shi. 2025. Hierarchical multi-label text classification of tourism resources using a label-aware dual graph attention network. *Information Processing & Management* 62, 1 (2025), 103952.
- [7] Dan Cogan, Zu-En Su, Oded Kenneth, and David Gershoni. 2023. Deterministic generation of indistinguishable photons in a cluster state. *Nature Photonics* 17, 4 (2023), 324–329.
- [8] Jiequan Cui, Beier Zhu, Qingshan Xu, Zhuotao Tian, Xiaojuan Qi, Bei Yu, Hanwang Zhang, and Richang Hong. 2025. Generalized Kullback-Leibler Divergence Loss. *arXiv preprint arXiv:2503.08038* (2025).
- [9] P. Ekman. 1992. Are there basic emotions? *Psychological Review* 99, 3 (1992), 550–553.
- [10] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [11] Jorge Luis Guerra, Carlos Catania, and Eduardo Veas. 2022. Datasets are not enough: Challenges in labeling network traffic. *Computers & Security* 120 (2022), 102810.
- [12] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 18234–18242.
- [13] Shadi Jaradat, Richi Nayak, Alexander Paz, Huthaifa I Ashqar, and Mohammad Elhenawy. 2024. Multitask learning for crash analysis: A fine-tuned llm framework using twitter data. *Smart Cities* 7, 5 (2024), 2422–2465.
- [14] Jamin Rahman Jim, Md Apon Riaz Talukder, Partha Malakar, Md Mohsin Kabir, Kamruddin Nur, and Mohammed Firoz Mridha. 2024. Recent advancements and challenges of NLP-based sentiment analysis: A state-of-the-art review. *Natural Language Processing Journal* (2024), 100059.
- [15] Francois Leonardi, Patrick Feldman, Matthew Almeida, William Moretti, and Charles Iverson. 2024. Contextual feature drift in large language models: An examination of adaptive retention across sequential inputs. (2024).
- [16] Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. 2023. Uniformer: Unifying convolution and self-attention for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 10 (2023), 12581–12600.
- [17] Yuming Li, Johnny Chan, Gabrielle Peko, and David Sundaram. 2024. An explanation framework and method for AI-based text emotion analysis and visualisation. *Decision Support Systems* 178 (2024), 114121.
- [18] Yang Li, Amirmohammad Kazemeini, Yash Mehta, and Erik Cambria. 2022. Multitask learning for emotion and personality traits detection. *Neurocomputing* 493 (2022), 340–350.
- [19] Zheng Li, Dawei Zhu, Qilong Ma, Weimin Xiong, and Sujian Li. 2024. EERPD: Leveraging Emotion and Emotion Regulation for Improving Personality Detection. *arXiv preprint arXiv:2406.16079* (2024).
- [20] Aixian Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).
- [21] Wenjie Liu, Zhijie Ren, and Liang Chen. 2025. Knowledge reasoning based on graph neural networks with multi-layer top-p message passing and sparse negative sampling. *Knowledge-Based Systems* (2025), 113063.
- [22] Wenjuan Liu, Zhengyan Sun, Subo Wei, Shunxiang Zhang, Guangli Zhu, and Lei Chen. 2024. PS-GCN: Psycholinguistic graph and sentiment semantic fused graph convolutional networks for personality detection. *Connection Science* 36, 1 (2024), 2295820.
- [23] Qing Luo, Wei Zeng, Manni Chen, Gang Peng, Xiaofeng Yuan, and Qiang Yin. 2023. Self-Attention and Transformers: Driving the Evolution of Large Language Models. In *2023 IEEE 6th International Conference on Electronic Information and Communication Technology (ICEICT)*. IEEE, 401–405.
- [24] Veronica Lynn, Niranjan Balasubramanian, and H Andrew Schwartz. 2020. Hierarchical modeling for user personality prediction: The role of message-level attention. In *Proceedings of the 58th annual meeting of the association for computational linguistics*. 5306–5316.
- [25] Walter Mischel and Yuichi Shoda. 1995. A cognitive-affective system theory of personality: reconceptualizing situations, dispositions, dynamics, and invariance in personality structure. *Psychological review* 102, 2 (1995), 246.
- [26] Max Murphy. 2024. Artificial Intelligence and Personality: Large Language Models' Ability to Predict Personality Type. *Emerging Media* (2024), 27523543241257291.
- [27] Benjamin Nelson, Ari Winbush, Steven Siddals, John Torous, Nick Allen, and Matthew Flathers. 2025. Evaluating the Performance of Large Language Models in Identifying Human Facial Emotions: GPT 4o, Gemini 2.0 Experimental, and Claude 3.5 Sonnet. (2025).
- [28] Ruthie Pliskin, Anat Ruhrman, and Eran Halperin. 2020. Proposing a multi-dimensional, context-sensitive approach to the study of ideological (a) symmetry in emotion. *Current Opinion in Behavioral Sciences* 34 (2020), 75–80.
- [29] J Prasanthi and G Anuradha. 2021. SURVEY ON PERSONALITY DETECTION USING DEEP LEARNING TECHNIQUES. In *2021 6th International Conference on Communication and Electronics Systems (ICCES)*. IEEE, 1–8.
- [30] Mustafa Safdari, Greg Serapio-Garcia, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matrić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184* (2023).
- [31] Aditya G Shanmukha, RS Shamyuktha, S Karan, Deepa Gupta, and Sujia Palaniswamy. 2024. Advancing Personality Detection through Word Embeddings and Deep Learning: An Examination Using the MBTI Dataset. In *2024 IEEE Recent Advances in Intelligent Computational Systems (RAICS)*. IEEE, 1–6.
- [32] Lingzhi Shen, Yunfei Long, Xiaohao Cai, Guanming Chen, Imran Razzak, and Shoaib Jameel. 2025. Less but Better: Parameter-Efficient Fine-Tuning of Large Language Models for Personality Detection. *arXiv preprint arXiv:2504.05411* (2025).
- [33] Lingzhi Shen, Yunfei Long, Xiaohao Cai, Guanming Chen, Yuhang Wang, Imran Razzak, and Shoaib Jameel. 2025. L4g: Self-supervised dynamic optimization for graph-based personality detection. *arXiv preprint arXiv:2504.02146* (2025).
- [34] Lingzhi Shen, Yunfei Long, Xiaohao Cai, Imran Razzak, Guanming Chen, Kang Liu, and Shoaib Jameel. 2025. Gamed: Knowledge adaptive multi-experts decoupling for multimodal fake news detection. In *Proceedings of the Eighteenth ACM International Conference on Web Search and Data Mining*. 586–595.
- [35] Zhouhao Sun, Xiao Ding, Li Du, Yunpeng Xu, Yixuan Ma, Yang Zhao, Bing Qin, and Ting Liu. 2025. Information Gain-Guided Causal Intervention for Autonomous Debiasing Large Language Models. *arXiv preprint arXiv:2504.12898* (2025).
- [36] Hossein Dabiryan Tehrani, Sara Yamini, and Alexander T Vazsonyi. 2024. Parenting styles and Big Five personality traits among adolescents: A meta-analysis. *Personality and Individual Differences* 216 (2024), 112421.
- [37] Teng Teng, Huifang Li, Yulin Fang, and Lingzhi Shen. 2022. Understanding the differential effectiveness of marketer versus user-generated advertisements in closed social networking sites: An empirical study of WeChat. *Internet Research* 32, 6 (2022), 1910–1929.
- [38] Iris Van Rooij, Olivia Guest, Federico Adolphi, Ronald de Haan, Antonina Kolokolova, and Patricia Rich. 2024. Reclaiming AI as a theoretical tool for cognitive science. *Computational Brain & Behavior* 7, 4 (2024), 616–636.
- [39] Di Wang, Ronghao Yang, Hanhu Liu, Haiqing He, Junxiang Tan, Shaoda Li, Yichun Qiao, Kangqi Tang, and Xiao Wang. 2022. HFENet: hierarchical feature extraction network for accurate landcover classification. *Remote Sensing* 14, 17 (2022), 4244.
- [40] Hongyu Wang, Dandan Zhang, Jun Feng, Lucia Cascone, Michele Nappi, and Shaohua Wan. 2024. A multi-objective segmentation method for chest X-rays based on collaborative learning from multiple partially annotated datasets. *Information Fusion* 102 (2024), 102016.
- [41] Y. Wang, D. Li, K. Funakoshi, and M. Okumura. 2023. Emp: Emotion-guided multi-modal fusion and contrastive learning for personality traits recognition. In *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval*. 243–252.

- [42] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [43] P William, N Yogeesh, Vishal M Tidake, Snehal Sumit Gondkar, K Vengatesan, et al. 2023. Framework for implementation of personality inventory model on natural language processing with personality traits analysis. In *2023 International Conference on Intelligent Data Communication Technologies and Internet of Things (IDCIoT)*. IEEE, 625–628.
- [44] Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia Sam Chao, and Derek Fai Wong. 2025. A survey on LLM-generated text detection: Necessity, methods, and future directions. *Computational Linguistics* (2025), 1–66.
- [45] Xuansheng Wu, Jiayi Yuan, Wenlin Yao, Xiaoming Zhai, and Ninghao Liu. 2025. Interpreting and steering llms with mutual information-based explanations on sparse autoencoders. *arXiv preprint arXiv:2502.15576* (2025).
- [46] Di Xue, Lifa Wu, Zheng Hong, Shize Guo, Liang Gao, Zhiyong Wu, Xiaofeng Zhong, and Jianshan Sun. 2018. Deep learning-based personality recognition from text posts of online social networks. *Applied Intelligence* 48, 11 (2018), 4232–4246.
- [47] Feifan Yang, Xiaojun Quan, Yunyi Yang, and Jianxing Yu. 2021. Multi-document transformer for personality detection. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 14221–14229.
- [48] Feifan Yang, Tao Yang, Xiaojun Quan, and Qinliang Su. 2021. Learning to answer psychological questionnaire for personality detection. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. 1131–1142.
- [49] Tao Yang, Jinghao Deng, Xiaojun Quan, and Qifan Wang. 2023. Orders are unwanted: dynamic deep graph convolutional network for personality detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 13896–13904.
- [50] Tao Yang, Tianyuan Shi, Fanqi Wan, Xiaojun Quan, Qifan Wang, Bingzhe Wu, and Jiaxiang Wu. 2023. PsyCoT: psychological questionnaire as powerful chain-of-thought for personality detection. *arXiv preprint arXiv:2310.20256* (2023).
- [51] Tao Yang, Feifan Yang, Haolan Ouyang, and Xiaojun Quan. 2021. Psycholinguistic tripartite graph network for personality detection. *arXiv preprint arXiv:2106.04963* (2021).
- [52] Hee Jun Yoon, Brent W Roberts, Madison N Sewell, Christopher M Napolitano, Christopher J Soto, Dana Murano, and Alex Casillas. 2024. Examining SEB skills’ incremental validity over personality traits in predicting academic achievement. *Plos one* 19, 1 (2024), e0296484.
- [53] Eric Zelikman, YH Wu, Jesse Mu, and Noah D Goodman. 2024. STaR: Self-taught reasoner bootstrapping reasoning with reasoning. In *Proc. the 36th International Conference on Neural Information Processing Systems*, Vol. 1126.
- [54] Haohao Zhu, Xiaokun Zhang, Junyu Lu, Youlin Wu, Zewen Bai, Changrong Min, Liang Yang, Bo Xu, Dongyu Zhang, and Hongfei Lin. 2024. Enhancing Textual Personality Detection toward Social Media: Integrating Long-term and Short-term Perspectives. *arXiv preprint arXiv:2404.15067* (2024).
- [55] Haohao Zhu, Xiaokun Zhang, Junyu Lu, Liang Yang, and Hongfei Lin. 2024. Integrating multi-view analysis: Multi-view mixture-of-expert for textual personality detection. In *CCF International Conference on Natural Language Processing and Chinese Computing*. Springer, 359–371.
- [56] Jianfeng Zhu, Ruoming Jin, and Karin G Coifman. 2025. Investigating Large Language Models in Inferring Personality Traits from User Conversations. *arXiv preprint arXiv:2501.07532* (2025).
- [57] Yangfu Zhu, Linmei Hu, Xinkai Ge, Wanrong Peng, and Bin Wu. 2022. Contrastive Graph Transformer Network for Personality Detection.. In *IJCAI*. 4559–4565.
- [58] Yangfu Zhu, Linmei Hu, Nianwen Ning, Wei Zhang, and Bin Wu. 2022. A lexical psycholinguistic knowledge-guided graph neural network for interpretable personality detection. *Knowledge-Based Systems* 249 (2022), 108952.
- [59] Nikola Zubić, Federico Soldá, Aurelio Sulser, and Davide Scaramuzza. 2024. Limits of Deep Learning: Sequence Modeling through the Lens of Complexity Theory. *arXiv preprint arXiv:2405.16674* (2024).