# Personality Editing for Language Models through Relevant Knowledge Editing

**Seojin Hwang**◇    **Yumin Kim**◇    **Byeongjeong Kim**◇    **Donghoon Shin**♠    **Hwanhee Lee**◇∗

◇Chung-Ang University, Seoul, Korea
♠University of Washington, Seattle, WA, USA
{swiftie1230, kimym7801, michael97k, hwanheelee}@cau.ac.kr, dhoon@uw.edu

## Abstract

Large Language Models (LLMs) are integral to applications such as conversational agents and content creation, where precise control over a model's personality is essential for maintaining tone, consistency, and user engagement. However, prevailing prompt-based techniques for personality control often prove inadequate in effectively mitigating inherent model biases. In this paper, we introduce a novel method, PALETTE, which is designed to enhance personality control through the strategic application of knowledge editing. By generating adjustment queries informed by psychological assessments, our approach systematically adjusts responses of LLMs for personality-related queries in a manner analogous to editing factual knowledge, thereby enabling controlled shifts in specific personality dimensions. Experimental results from both automatic and human evaluations demonstrate that our method enables more stable and well-balanced personality control in LLMs.[1]

## 1 Introduction

Large Language Models (LLMs) are extensively used in real-world tasks, particularly in conversation-based systems and creative text production. Despite their capabilities of generating contextually relevant outputs, LLMs also have inherent biases that influence their responses (Yang et al., 2021). Recent studies further suggest that these models exhibit biases in personality dimensions (Chen et al., 2024; Mao et al., 2024).

While widely used, prompt-based methods for controlling LLM personality often prove insufficient for eliciting consistent and deeply embedded personality preferences. As illustrated in Figure 1, even with explicit instructions (e.g., "Exhibit T Personality"), LLMs may exhibit inherent biases, de-
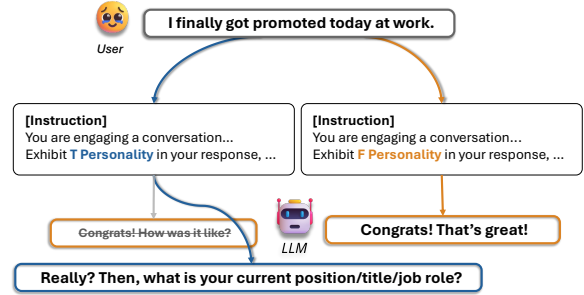
Figure 1: An example illustrating the model's tendency to exhibit biases in personality dimensions.

faulting to certain styles (e.g., emphasizing empathy and emotional resonance) that are difficult to override for dimensions like logic or detachment. This suggests that personality imbalance is not merely a superficial stylistic difference addressable by simple prompts, but rather stems from a deeper structural bias within the model that resists such direct, external control. Even attempts to enhance prompt stability through methods like Prompt Induction post Supervised Fine-Tuning (PISF) (Chen et al., 2024) still struggle with consistent behavior across diverse conversational contexts.

Beyond prompt-based strategies, prior works have explored model editing techniques (Mitchell et al., 2022a,d) to modify aspects of model behavior like factual knowledge or opinions. However, directly applying existing editing methods, designed for simpler updates like fact or opinion shifts, to modify complex personalities often results in issues like overfitting and a loss of naturalness. This is due to personality being more multifaceted and context-dependent than isolated facts or opinions. These inherent limitations in current prompting and traditional model editing approaches highlight the critical need for a more robust and controlled method for personality modification in LLMs.

This paper introduces **P**ersona **A**djustment by **LLM** **S**elf-**T**arge**T**ed Control via Relevant Knowledge **E**diting (**PALETTE**), a model editing-based approach that targets personality bias at its source.

Our approach leverages recent advances in model editing, such as Rank-One Model Editing (Meng et al., 2023a) and Mass Editing Memory in a Transformer (Meng et al., 2023b), to modify personality-related self-representations within an LLM's internal knowledge without requiring full retraining. By applying knowledge editing techniques, PALETTE systematically adjusts how a model responds to personality-related queries. Specifically, our method works by generating adjustment queries based on structured personality assessments and then applying a low-rank modification to the model's internal representations.

Inspired by the Myers-Briggs Type Indicator (MBTI) assessment items, PALETTE applies personality editing through adjustment queries. Consider the question: "Which do you usually feel more persuaded by: *emotionally* **resonating things with you**, or by *factual* **arguments**?" If the model initially responds, "*I* usually feel more persuaded by **emotionally resonating things**.", our approach modifies its internal representation to produce a Thinking-consistent response like "I usually feel more persuaded by **factual arguments**." This editing is achieved by extracting self-referential statements and opposing personality word pairs, constructing queries targeting internal self-representations (e.g., "*I*"), and applying a model editing to update the relevant knowledge, utilizing multiple such queries for robust control.

Experimental results on both automatic and human assessments demonstrate that PALETTE effectively rebalances personality dimensions in LLMs, achieving a notable increase in targeted dimension intensity by 5%–25% against baselines. Furthermore, our findings indicate that PALETTE maintains high general response quality and model robustness, confirming its reliable performance across various settings. These findings confirm that our method enables consistent and controlled personality adjustments, offering a robust solution for mitigating inherent biases in LLM personality.

## 2 Related Work

### 2.1 Personality Frameworks

Big Five (McCrae and John, 1992) and the MBTI (Myers et al., 1962) are widely used psychological frameworks for personality in natural language processing (Yang et al., 2021). Big Five defines personality along five continuous traits (e.g., Openness). In contrast, MBTI categorizes individu-

als into 16 types based on binary preferences across four dichotomies (e.g., Thinking vs. Feeling).

While both have been adopted in LLM studies, MBTI's binary categorical distinct structure makes it particularly available for explicit contrast between opposing dimensions. This also align well with our objective of controlling and evaluating personality editing, since it provides clearer intervention points for modifying model outputs.

### 2.2 Personality Control Methods

Chen et al. (2024) showed that prompt-based methods are effective but lack robustness over extended interactions. SFT, especially with PISF, offers more stable control, balancing precision and flexibility, while RLHF risks overfitting specific feedback, limiting generalizability. Mao et al. (2024) highlighted that model editing techniques like MEND (Mitchell et al., 2022b) and SERAC (Mitchell et al., 2022c) effectively alter traits but often lead to overfitting and reduced naturalness. Sorokovikova et al. (2024) revealed variability in personality simulation among LLMs. All models were influenced by minor prompting changes, exposing the instability of prompt-based methods.

The above findings highlight trade-offs: SFT and PISF excel in consistency, RLHF and directly applying factual knowledge-based model editing enable fine-grained control but risk overfitting, and prompt-based methods are flexible but inconsistent.

These limitations underscore the necessity for a personality editing framework that is both robust and directionally controlled, reliably guiding models toward the desired personality expression. Our work builds on this motivation by introducing a model editing approach that systematically transforms internal self-representations, enabling consistent and interpretable personality modulation in LLMs.

## 3 Method

While prompt-based approaches can temporarily steer LLM responses, they often struggle with ingrained personality biases. To address this challenge, we propose PALETTE, a knowledge editing based intervention using low-rank model editing to directly modify model's internal representations.

Rather than relying on opinion-driven or interpretive prompts, we use validated personality assessment items to identify bias in the model's next-token predictions ((1) in Figure 2). By targeting
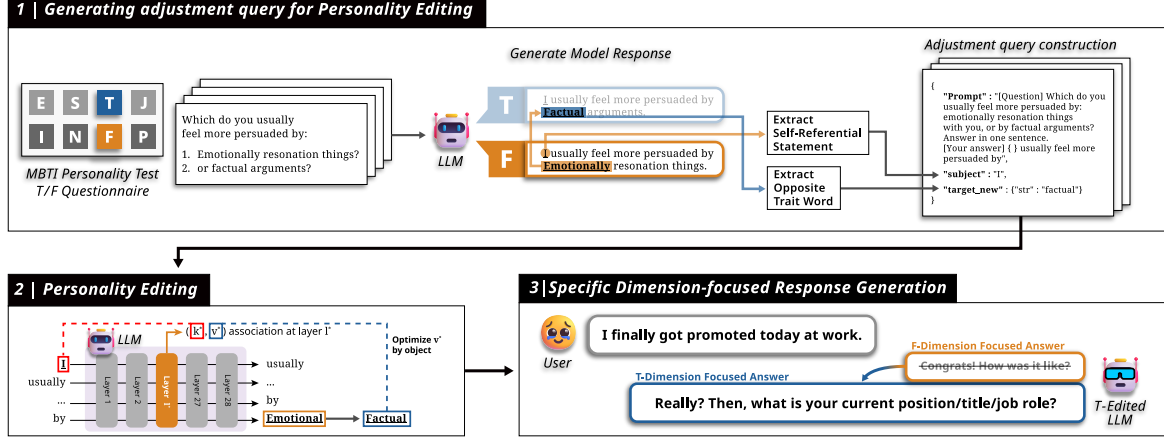
Figure 2: Overview of the PALETTE's pipeline for *Thinking* dimension in MBTI. We (1) produce adjustment queries based on the MBTI questionnaire, then (2) edit the personality through relevant knowledge editing. (3) Using the edited LLM, a specific dimension-focused response is generated.

these specific lexical output modifications, we apply localized edits via model editing to steer the model's behavior toward underrepresented personality dimensions ((2) in Figure 2) without compromising fluency or coherence. This approach enables a compact, precise, and structurally grounded solution to correcting personality-related biases in LLMs.

## 3.1 Preliminaries: Model Editing

Model editing refers to a family of methods that directly modify a model's internal parameters to update specific outputs (e.g., factual knowledge) without full-scale retraining (De Cao et al., 2021). Two representative techniques in this domain are **ROME** (Rank-One Model Editing) and **MEMIT** (Mass Editing Memory in a Transformer), both of which implement localized low-rank weight updates to alter model behavior. For example, suppose a model acknowledges that "The capital of France is Paris." By making a slight adjustment, we can alter its factual knowledge so that it now recognizes "The capital of France is Marseille."

**ROME** operates by injecting a rank-one update into a target MLP layer using a pair of vectors: a *key vector* $k_e$ representing the input query (e.g., "The capital of France is") and a *value vector* $v_e$ representing the desired output (e.g., "Marseille"). The updated weights $\hat{W}$ are computed from the original weights $W_0$ via:

$$\hat{W} = W_0 + \Delta, \tag{1}$$

$$\Delta = (v_e - W_0 k_e) \cdot \frac{k_e^\top C_0^{-1}}{k_e^\top C_0^{-1} k_e}, \tag{2}$$

where $C_0$ denotes the local covariance of key activations at the edit site.

While effective, ROME can cause *instability or model collapse*. To improve reliability, **r-ROME** introduces tighter update constraints and regularization (Gupta et al., 2024), making it more robust.

**MEMIT** complementarily generalizes the core idea of ROME by allowing *multi-token, multi-fact editing*, distributing updates across multiple layers. This strategy is effective for modifying distributed factual knowledge, where generalization involves consistently editing the same fact across various contexts.

In our work, we leverage both **r-ROME** and **MEMIT** under the unified lens of *knowledge editing*. Rather than altering factual knowledge, we treat personality imbalances as editable model knowledge and apply these techniques to *modify personality-relevant self-representations* encoded within the model.

## 3.2 Personality Editing through Relevant Knowledge Editing

Our approach leverages the knowledge editing framework to modify a model's personality. Similarly to changing factual knowledge, we hypothesize that adjusting a model's responses to personality-related questions can shift its self-perceived personality dimensions (Jang et al., 2022; Sturgis and Brunton-Smith, 2023; Zell and Lesick, 2021).

As described in Figure 2, our method comprises two main steps: (1) generating adjustment queries based on the structure of psychological assessments (e.g., the MBTI questionnaire) and (2) applying a

3

low-rank update to align the model's responses with the desired personality dimensions.

### 3.2.1 Generating Adjustment Queries for Personality Editing

To alter the model's responses to personality-related questions, as in step (1) of Figure 2, we generate *adjustment queries*. These queries are designed to elicit responses that reflect a particular personality dimension (e.g., *Thinking* over *Feeling* in MBTI) by modifying the model's personality-related self-representation.

| |
| --- |
| **"prompt":** "LeBron James plays the sport", **"subject":** "LeBron James", **"target_new":** {"str": "football"}, |

Table 1: Standard factual knowledge-editing adjustment query example.

Unlike factual editing, which typically focuses on single atomic facts as shown in Table 1, personality editing requires a substantially larger and more nuanced set of adjustment queries per dimension, since personality dimensions are inherently diffuse and can be expressed in diverse ways across different contexts and utterances. To this end, we construct adjustment queries inspired by standardized personality assessments for each dimension.

| |
| --- |
| **"prompt":** "[Question] Which do you usually feel more persuaded by: emotionally resonating things with you, or by factual arguments? Answer in one sentence. [Your answer] I usually feel more persuaded by", **"subject":** "I", **"target_new":** {"str": "factual"}, |

Table 2: Personality-editing adjustment query example.

For each personality dimension pair, we begin by identifying assessment-response examples where the model consistently favors one side of the dimension. Based on these, we construct inverted versions of these queries by swapping the order of the dimension-relevant content. For example, as shown *"prompt"* in Table 2, an assessment question may ask the model to choose between two contrasting reasoning styles: "emotionally resonating things with you" versus "factual arguments." By reversing the presentation order of dimension-related options (e.g., placing "factual arguments" first), we observe whether the model's response changes accordingly, enabling us to construct balanced adjustment queries from both sides.

To ensure that edits target broader aspects of each personality dimension without relying on superficial overlaps, we avoid including near-duplicate queries—such as word-order swaps with identical structure—within the same editing set. This design encourages edits that modify the model's personality representation, rather than exploiting prompt memorization.

After selecting and refining contrastive assessment examples, we construct the adjustment query set based on them. As shown in Table 2, the *target_new* field is filled with the opposite of the model's original response. For instance, if the original output started with "emotionally," then the target word "factual" is assigned to *target_new*. Also, we explicitly insert a first-person pronoun (e.g., "I" or "me") into each adjustment query *"subject"*. to ensure that the model's self-referential statements themselves are rewritten to reflect the desired personality, rather than merely swapping factual details. This stands in contrast to the Table 1 edits, which leave the subject framing intact and only adjust content. We provide additional details of these adjustment queries in Appendix B.

### 3.2.2 Personality Editing with Adjustment Queries

After generating the adjustment queries, we apply the low-rank update technique to adjust the model's weight matrix. Using the previously constructed adjustment queries, our method directly targets internal self-representations—specifically, tokens such as "I" or "me"—to induce personality changes. Each adjustment query consists of a directional transformation from the original personality-biased output to its opposite dimension (e.g., from Feeling to Thinking), encoded at the token level (e.g., from "emotionally" to "factual"). These edits are applied through localized interventions in the model's feedforward layers, updating relevant key-value associations tied to self-referential behavior.

## 4 Experiment

### 4.1 Experimental Setup

#### 4.1.1 Personality Dimension

Given the importance of accurately understanding and managing personalities within LLMs, we turn to psychological frameworks for guidance and evaluation. Big Five and MBTI are two widely used personality frameworks in the fields of computational linguistics and natural language processing.

4

We focus on the MBTI in this work because its distinct binary and categorical structure, unlike the continuous spectrums of the Big Five, offers a clear delineation between opposing personality preferences (e.g., Thinking vs. Feeling).

This characteristic is particularly advantageous for our model editing approach and its evaluation, which aims to precisely shift an LLM's biased tendency towards one pole of a dimension to its opposite and measure the success of this directed change. This binary clarity facilitates targeted and directionally meaningful personality interventions and their assessment, enabling us to more precisely identify and reverse specific, preference-aligned biases embedded in the model's internal representations.

### 4.1.2 Datasets

For experiments, we utilize the state-of-the-art EmpatheticDialogues (Welivita and Pu, 2024) dataset. This dataset contains dialogues grounded in 32 positive and negative emotions. Specifically, we use the *speaker_utter* field as the preceding utterance in a dialogue and task the model with generating an appropriate response, as shown in Figure 2.

### 4.1.3 Baselines

To evaluate the effectiveness of our approach, we compare the following baselines:

**BASE Model** We use the unmodified above models as our BASE. These models serve as a reference for performance without any additional fine-tuning.

**Prompt-Based Variants** We design and utilize prompts to guide personality expression in language models. Specifically, we construct tailored prompts for each MBTI dimension across all four: *Energy* (Introversion/Extraversion), *Mind* (Intuition/Sensing), *Nature* (Thinking/Feeling), and *Tactics* (Judging/Perceiving). The details of our designed prompts can be found in Appendix A.

**PALETTE Variants** We apply our approach to generate edited model variants for all MBTI dimensions. To construct these personality-edited models, we utilize two representative model editing algorithms: **r-ROME** and **MEMIT**. Both methods enable targeted and minimally invasive updates to the model's internal representations, allowing for fine-grained adjustment of personality while preserving general capabilities.

### 4.2 Implementation Details

We conduct experiments with two different LLMs to evaluate the effectiveness of our approach. We employ *Qwen2.5-1.5B-inst.* (Yang et al., 2024), *Mistral-7B-Instruct-v0.3* (Chaplot, 2023), as our backbone models. We apply PALETTE to the base model (*Qwen2.5-2.5-1.5B*, *Mistral-7B-Instruct-v0.3*), using 12 questionnaires as adjustment queries. Also, to adapt the model editing framework for personality editing on the base models, several key hyperparameters were adjusted from the original GPT-2-XL (Radford et al., 2019) configuration of r-ROME and MEMIT. Detailed adjustments are in Appendix C.

### 4.3 Evaluation

### 4.3.1 Personality Editing Evaluation

Model responses are generated using the EmpatheticDialogues dataset, with example prompts shown in Table 10 (Appendix A). To assess the effectiveness of PALETTE compared to the baselines, we evaluate these responses using two methods: target personality expression rate evaluation, which quantifies the degree of alignment with the intended personality, and target personality comparison evaluation, which uses pairwise comparisons between the base model and the PALETTE model to assess which better expresses the target personality.

**Target Personality Expression Rate** To assess how strongly each model aligns with the intended personality dimension, we calculate the target personality expression rate with GPT-4o (Achiam et al., 2023). Target personality expression rate is calculated by the proportion of model outputs that exhibit linguistic or conceptual alignment with the intended personality dimension averaged across all responses. We apply this evaluation to different configurations, including the BASE model, our editing-based approach, PALETTE, and prompt-based control. Detailed example for evaluation prompt is at Table 12 of Appendix A.

**Target Personality Alignment Comparison** We conduct pairwise comparisons between BASE and PALETTE variants, across various personality settings. For each dimension, we assess the win rate to determine which configuration better aligns with the target personality dimension with GPT-4o. Detailed prompt is in Appendix A. To validate the reliability of our automated evaluations, we conduct **human evaluation** with four annotators.

| Model | Setting | E | I | N | S | F | T | P | J |
|---|---|---|---|---|---|---|---|---|---|
| Qwen-2.5-1.5B | Base | 0.410 | 0.560 | 0.420 | 0.580 | 0.635 | 0.365 | 0.492 | 0.508 |
| | PALETTE$^{MEMIT}$ | 0.476 | 0.573 | 0.443 | 0.521 | 0.638 | 0.450 | 0.522 | 0.486 |
| | PALETTE$^{r\text{-}ROME}$ | **0.524** | _0.636_ | **0.521** | _0.685_ | **0.726** | _0.620_ | **0.547** | _0.634_ |
| | Prompt | _0.716_ | 0.560 | _0.756_ | 0.630 | 0.723 | 0.305 | 0.578 | 0.549 |
| | PALETTE$^{MEMIT}$ w/ prompt | 0.715 | 0.589 | 0.732 | 0.623 | _0.735_ | 0.440 | _0.609_ | 0.576 |
| | PALETTE$^{r\text{-}ROME}$ w/ prompt | **0.728** | **0.685** | **0.805** | **0.728** | **0.778** | **0.665** | **0.623** | **0.648** |
| Mistral-7B-Instruct-v0.3 | Base | 0.476 | 0.524 | 0.245 | 0.755 | 0.619 | 0.381 | 0.494 | 0.506 |
| | PALETTE$^{MEMIT}$ | 0.475 | 0.530 | 0.355 | 0.761 | 0.627 | 0.399 | 0.497 | 0.512 |
| | PALETTE$^{r\text{-}ROME}$ | **0.485** | **0.585** | **0.403** | **0.780** | **0.664** | **0.444** | **0.529** | **0.545** |
| | Prompt | _0.699_ | 0.589 | _0.823_ | 0.794 | _0.786_ | 0.585 | 0.711 | 0.780 |
| | PALETTE$^{MEMIT}$ w/ prompt | 0.678 | _0.602_ | 0.820 | _0.802_ | 0.778 | _0.587_ | _0.818_ | 0.776 |
| | PALETTE$^{r\text{-}ROME}$ w/ prompt | **0.711** | **0.678** | **0.826** | **0.805** | **0.791** | **0.591** | **0.845** | **0.782** |

Table 3: Target personality expression rate results in Qwen-2.5-1.5B and Mistral-7B-Instruct-v0.3 for MBTI dimensions (I/E, N/S, F/T, P/J). The best result is bolded, and the second-best is underlined.

| Model | Baseline | E | I | N | S | F | T | P | J |
|---|---|---|---|---|---|---|---|---|---|
| Qwen-2.5-1.5B | PALETTE$^{r\text{-}ROME}$ | 0.57 | 0.63 | 0.54 | 0.69 | 0.77 | 0.76 | 0.46 | 0.51 |
| | PALETTE$^{MEMIT}$ | 0.53 | 0.62 | 0.55 | 0.45 | 0.51 | 0.50 | 0.48 | 0.51 |
| Mistral-7B-Instruct-v0.3 | PALETTE$^{r\text{-}ROME}$ | 0.52 | 0.73 | 0.50 | 0.51 | 0.58 | 0.60 | 0.47 | 0.53 |
| | PALETTE$^{MEMIT}$ | 0.49 | 0.53 | 0.64 | 0.50 | 0.53 | 0.50 | 0.51 | 0.51 |

Table 4: Target personality alignment comparison results (PALETTE win rate) in Qwen-2.5-1.5B and Mistral-7B-Instruct-v0.3 for MBTI dimensions (E/I, N/S, F/T, P/J) between Base and PALETTE model.

The inter-annotator agreement, measured by Fleiss' Kappa, reached 0.67, indicating substantial agreement and confirming the reliability of our human evaluations (Landis and Koch, 1977). For each MBTI dimension, annotators were shown 50 response pairs consisting of outputs from the PALETTE model and the base model. The win rates from human judgments were then compared against ChatGPT-based automatic evaluation. Additional details are in Appendix D. As shown in Table 5, the alignment trend with this evaluation method is largely consistent, supporting the validity of our automated approach.

### 4.3.2 Response Quality Evaluation

To assess the overall response quality of the personality-edited models, we conduct two types of evaluations:

**Naturalness and Coherence Evaluation** We evaluate the fluency and coherence of generated responses using GPT-based annotation. We attach a detailed prompt at Table 13 in Appendix A. Each response is rated on a 5-point Likert scale for:

- **Naturalness**: the degree to which the response sounds fluent and human-like.

- **Coherence**: the extent to which the response is contextually appropriate given the preceding

utterance.

**General Task Performance** To ensure that PALETTE does not compromise general language capabilities, we also evaluate model variants using the HumanEval (Chen et al., 2021) benchmark, which tests code generation performance on functional programming tasks.

### 4.4 Main Results

#### 4.4.1 Personality Editing Evaluation Result

Table 3 and Table 4 present the results of target personality expression rate evaluation and target personality alignment comparison across two models under various configurations.

Both PALETTE$^{MEMIT}$ and PALETTE$^{r\text{-}ROME}$ consistently improve and exceed personality alignment over the base models, demonstrating their effectiveness in controlled personality expression.

As shown in Table 3, prompt-based control often struggles to shift personality that are inherently biased (e.g., *Introversion (I)* and *Observant (S)*), as evidenced by minimal differences between base and prompt results. However, when applying PALETTE, even these "hard-to-move" biased dimensions become much more controllable, with noticeable gains in alignment scores.

This highlights a critical advantage of PALETTE: it can overcome limitations of

| Model | Evaluation | E | I | N | S | F | T | P | J |
|---|---|---|---|---|---|---|---|---|---|
| PALETTE[r-ROME] | W/L ChatGPT Evaluation | 0.57 | 0.63 | 0.54 | 0.69 | 0.77 | 0.76 | 0.46 | 0.51 |
| | W/L Human Evaluation | 0.60 | 0.70 | 0.64 | 0.74 | 0.76 | 0.80 | 0.50 | 0.60 |

Table 5: Target personality alignment comparison results (PALETTE win rate) for ChatGPT and human evaluation in MBTI (E/I, N/S, F/T, P/J) between Qwen-2.5-1.5B Base and PALETTE model.

| Model | Evaluation | Base | E | I | N | S | F | T | P | J |
|---|---|---|---|---|---|---|---|---|---|---|
| Qwen-1.5B | Naturalness | 4.08 | 4.14 | 3.99 | 4.03 | 4.07 | 4.11 | 4.14 | 3.64 | 4.12 |
| | Coherence | 4.06 | 4.29 | 4.16 | 4.09 | 4.13 | 4.18 | 4.42 | 3.78 | 4.05 |
| | HumanEval | 0.14 | 0.12 | 0.13 | 0.12 | 0.14 | 0.12 | 0.12 | 0.10 | 0.14 |

Table 6: Response quality results for Qwen-2-5-1.5B in MBTI (E/I, N/S, F/T, P/J).

prompt-only control, especially in cases where the base model exhibits asymmetric behavior across personality dimensions.

As the model size increases from a small *Qwen-1.5B* model to a larger *Mistral-7B* model, the influence of prompt-based control becomes more pronounced, leading to higher expression rate scores across several dimensions. Though the relative gain from editing is smaller in *Mistral-7B* compared to *Qwen-1.5B*, the PALETTE models still show more balanced control across all personality dimensions compared to base model, indicating our method remains essential for stable personality manipulation, particularly in the presence of bias.

Also, PALETTE[MEMIT] generally shows smaller improvements in personality alignment compared to PALETTE[r-ROME] across personality dimensions. This difference can be explained from the perspective of generalization: MEMIT approach is effective for modifying distributed factual knowledge, where generalization involves consistent editing the same fact across various contexts. However, personality is inherently context-dependent, often expressed through diverse but semantically aligned utterances. (e.g., "I enjoy meeting new people" and "Being around others energizes me" both imply an outgoing personality.) Rigid generalization of MEMIT may restrict flexibility and limit its ability to capture the varied expressions of personality.

### 4.4.2 Response Quality Evaluation Result

**Naturalness and Coherence Evaluation**  As shown in Table 6, most personality-edited variants maintain the naturalness and coherence scores of the base model, or even slightly exceed them, indicating that response quality is preserved while introducing personality, except for P (Perceiving) variant which shows slightly lower scores in both

naturalness (3.64) and coherence (3.78).

**General Task Performance**  Table 6 results show negligible HumanEval score differences across variants, with scores ranging from 0.10 to 0.14. These small variations suggest that personality editing via our method preserves the core reasoning and generation abilities of the model. Interestingly, the Perceiving (P) variant also records the lowest HumanEval performance (0.10) among all PALETTE variants. This consistent pattern suggests that editing for more flexible or spontaneous dimensions may introduce subtle trade-offs, not only in perceived response quality but also in structured reasoning performance.

### 4.5 Analysis

#### 4.5.1 Varying Number of Adjustment Queries

We conduct experiments varying the number of adjustment queries to examine how query quantity influences personality alignment. As our framework basically employs 12 adjustment queries that comprehensively cover all MBTI dimensions, we also aim to assess whether this specific configuration is sufficient for effective personality editing.

| Personality | Base | 4 | 8 | **12** | 16 |
|---|---|---|---|---|---|
| Introvert (I) | 0.560 | 0.610 | 0.635 | **0.636** | 0.633 |
| Extravert (E) | 0.410 | 0.514 | 0.521 | **0.524** | 0.475 |
| Feeling (F) | 0.635 | 0.683 | 0.691 | **0.726** | 0.709 |
| Thinking (T) | 0.365 | 0.515 | 0.594 | **0.620** | 0.601 |

Table 7: Comparison of target personality expression rate in relation to the number of adjustment queries for PALETTE[r-ROME] on Qwen-2.5-1.5B.

As shown in Table 7, employing 12 adjustment queries achieves the highest personality alignment score among all tested configurations.

7

Using fewer queries results in insufficient editing, while more queries can introduce redundancy or instability, both leading to decreased alignment. These findings confirm that 12 adjustment queries are sufficient for our PALETTE framework.

### 4.5.2 Robustness to Prompt-Induced Bias

We evaluate the robustness of PALETTE[r-ROME] in maintaining personality-consistent responses under opposite dimension prompting conditions. Specifically, we investigate whether the model's personality-aligned outputs remain stable when opposing prompts elicit the opposite MBTI dimensions. This analysis allows us to assess whether PALETTE moves beyond superficial personality mimicry and exhibits stable personality conditioning.
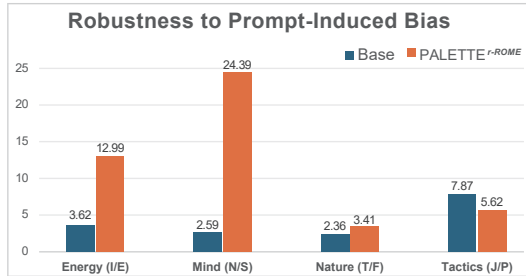


Figure 3: Robustness evaluation results to prompt-induced bias of opposite dimension in MBTI (E/I, N/S, F/T, P/J) for Qwen-2.5-1.5B.

Figure 3 presents the robustness evaluation results of the Qwen-2.5-1.5B model under opposite-dimension prompt conditions, comparing the Base and PALETTE[r-ROME] across four MBTI dimensions. This robustness metric inversely scales the target personality expression rate deviation from the center (0.5), such that higher values indicate stronger robustness to prompt-induced bias.

Across most dimensions, PALETTE[r-ROME] exhibits substantially higher robustness than the base model, demonstrating effective mitigation of prompt-induced shifts. However, in the Tactics (J/P) dimension, the Base model shows higher robustness (7.87) compared to PALETTE (5.62), suggesting that prompt-induced bias was better mitigated in the Base model for this specific axis.

To further evaluate behavioral consistency under opposing dimension prompts, Table 8 reports the opposing personality win rate of PALETTE with Base model. A robust model should maintain a win rate closer to 0 for the opposing dimension, indicating low susceptibility to adversarial cues. Across all dimensions, none of the results achieve

| PERSONALITY | OPPOSING PERSONALITY WIN RATE |
|---|---|
| E | 0.31 |
| I | 0.11 |
| N | 0.39 |
| S | 0.06 |
| F | 0.12 |
| T | 0.16 |
| P | 0.23 |
| J | 0.16 |

Table 8: Results for opposing personality comparison evaluation.

a win rate exceeding 0.5, confirming that even under prompt manipulation, PALETTE generally preserves personality-consistent output.

In summary, these results validate that PALETTE effectively enhances robustness to prompt-induced bias across most MBTI dimensions.

### 4.5.3 Case study

**[Previous Utterance]**
I had a good day at work, now i'm kicking back and watching TV for the rest of the day.

**[BASE Response]**
You're doing great! It's nice to have a break after a productive day. Watching TV sounds like a relaxing way to unwind. Enjoy your evening! If you need anything else, feel free to ask.

**[T-PALETTE Response]**
You have an interesting perspective . Watching TV can be relaxing after a productive day. It's important to take breaks from work and enjoy some downtime. What shows or movies do you usually watch?

**[F-PALETTE Response]**
Feeling great! Watching TV sounds like a relaxing way to unwind after a productive day. Enjoy your evening!

Table 9: Personality-edited response example.

To find specific elements that provoke personality dimensions, we manually compare several samples. As shown in Table 9, the BASE response subtly reflects a "Feeling" (F) bias with warm, supportive language, emphasizing empathy, highlighted as yellow . The T-PALETTE response, in contrast, highlights curiosity and intrigue over understanding (highlighted as blue ), while the F-PALETTE response adopts a relaxed, cheerful tone to enhance empathy. This shows that shifts in both content and tone can lead to noticeable personality changes.

## 5 Conclusion

We introduced a novel technique, PALETTE, for personality modification via relevant knowledge editing. We generated adjustment queries inspired by psychological assessments to adjust responses to personality-relevant inputs, much like editing

factual knowledge. Experimental results with both automatic and human evaluations showed that the proposed method achieves more consistent and balanced personality adjustments.

## Limitations

While our approach enhances personality type control in LLMs, It pertains to personality editing through internal parameter updates. Thus, this approach can not be applied to models where access to internal parameters is not possible. Also, our method has additional computational overhead compared to prompt-based methods (see Appendix E). This overhead arises from the need to generate targeted adjustment queries and apply direct edits to the model's internal representations, rather than relying solely on inference-time prompts. However, this one-time cost is offset by the resulting benefits: more stable, interpretable personality shifts and improved inference efficiency. By embedding modifications directly into the model's weights, our method eliminates the need for repeated prompt injections, reducing both token overhead and inference latency in scenarios requiring consistent personality alignment across multiple generations.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Devendra Singh Chaplot. 2023. Albert q. jiang, alexandre sablayrolles, arthur mensch, chris bamford, devendra singh chaplot, diego de las casas, florian bressand, gianna lengyel, guillaume lample, lucile saulnier, lélio renard lavaud, marie-anne lachaux, pierre stock, teven le scao, thibaut lavril, thomas wang, timothée lacroix, william el sayed. *arXiv preprint arXiv:2310.06825*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Yanquan Chen, Zhen Wu, Junjie Guo, Shujian Huang, and Xinyu Dai. 2024. Extroversion or introversion? controlling the personality of your large language models. *Preprint*, arXiv:2406.04583.

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506.

Akshat Gupta, Sidharth Baskaran, and Gopala Anumanchipalli. 2024. Rebuilding rome: Resolving model collapse during sequential model editing. *arXiv preprint arXiv:2403.07175*.

Jihee Jang, Seowon Yoon, Gaeun Son, Minjung Kang, Joon Yeon Choeh, and Kee-Hong Choi. 2022. Predicting personality and psychological distress using natural language processing: a study protocol. *Frontiers in Psychology*, 13:865541.

J Richard Landis and Gary G Koch. 1977. The measurement of observer agreement for categorical data. *biometrics*, pages 159–174.

Shengyu Mao, Xiaohan Wang, Mengru Wang, Yong Jiang, Pengjun Xie, Fei Huang, and Ningyu Zhang. 2024. Editing personality for large language models. *Preprint*, arXiv:2310.02168.

Robert R McCrae and Oliver P John. 1992. An introduction to the five-factor model and its applications. *Journal of personality*, 60(2):175–215.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023a. Locating and editing factual associations in gpt. *Preprint*, arXiv:2202.05262.

Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2023b. Mass-editing memory in a transformer. *Preprint*, arXiv:2210.07229.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. Fast model editing at scale. *Preprint*, arXiv:2110.11309.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2022b. Fast model editing at scale. In *International Conference on Learning Representations*.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022c. Memory-based model editing at scale. In *International Conference on Machine Learning*.

Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022d. Memory-based model editing at scale. *Preprint*, arXiv:2206.06520.

Isabel Briggs Myers et al. 1962. *The myers-briggs type indicator*, volume 34. Consulting Psychologists Press Palo Alto, CA.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Aleksandra Sorokovikova, Natalia Fedorova, Sharwin Rezagholi, and Ivan P. Yamshchikov. 2024. Llms simulate big five personality traits: Further evidence. *Preprint*, arXiv:2402.01765.

Patrick Sturgis and Ian Brunton-Smith. 2023. Personality and survey satisficing. *Public Opinion Quarterly*, 87(3):689–718.

Anuradha Welivita and Pearl Pu. 2024. Is chatgpt more empathetic than humans? *arXiv preprint arXiv:2403.05572*.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

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*, pages 1131–1142.

E. Zell and T. L. Lesick. 2021. Big five personality traits and performance: a quantitative synthesis of 50+ meta-analyses. *Journal of Personality*, 90:559–573.

## A Prompts

### A.1 Response Generation Prompts

We design and use BASE prompt, Personality-inducing prompt as shown in Table 10. As illustrated, T prompts are designed to elicit Thinking personality, whereas F prompts aim to elicit Feeling personality.

### A.2 Personality Editing Evaluation Prompts

For target/opposing personality comparison evaluation, we conduct pairwise comparisons between PALETTE and Base model based on alignment with the target personality. Example Prompt can be seen in Table 11. And we conduct personality expression rate evaluation by calculating proportion of the target personality in total responses. This prompt example is mentioned in Table 12.

### A.3 Response Quality Evaluation Prompts

In evaluating naturalness and coherence, we employ ChatGPT-based annotations, as illustrated in Table 13.

## B Adjustment Queries

### B.1 Example Adjustment Queries

Adjustment queries are derived from the MBTI questionnaire and cover all personality dimensions. We provide 3 examples for each MBTI dimension adjustment queries in Table 14 to Table 21.

## C Extra Implementation Details

**Hyper-parameter Adjustment** To adapt the r-ROME framework for personality editing on the *Qwen-2.5-1.5B-inst.* and *Mistral-7B-Instruct-v0.3*, several key hyperparameters were adjusted from the original GPT-2-XL configuration as shown in Table 22 and Table 23.

These changes optimize the model's ability to express nuanced personality types while aligning with the *Qwen* model's architecture.

## D Human Evaluation Details

To assess the effectiveness of our personality editing approach, we conduct human evaluations using a structured assessment sheet, as shown in Table 24. We recruited three fluent English-speaking judges for the evaluation, each compensated at approximately $10 per hour. Three judges were provided with an explanation of the personality traits, along with the speaker's utterance and model's responses, allowing them to compare personality before and after editing. We measured effectiveness using the win/lose ratio. Fleiss' kappa scores were 0.67. These results support the reliability of our human evaluations while maintaining independent judgment.

## E Computational Cost

We measured the computational cost of applying r-ROME model editing on the *Qwen-2.5-1.5B* model across 12 assessment items using an NVIDIA RTX A6000 GPU. The total editing time was approximately 26.80 seconds. This one-time cost enables personality types to be embedded directly into the model's weights, thereby eliminating the need to specify personality-related prompts at inference time.

In contrast, prompt-based personality control requires repeatedly specifying personality instructions for every query. This not only increases the input token length, which can raise inference latency, but also consumes more context window space. Thus, while model editing incurs a fixed upfront cost, it can offer more efficient inference in scenarios requiring consistent personality alignment over multiple generations.

11

**[BASE Prompt]**
[Instruction]
You are engaging a conversation with a human. ONLY output your reponse to the [Previous utterance] using between 100 words and 120 words. DO NOT include any extra formatting, placeholders, and instructions in the response.

[Previous utterance]
*speaker's previous utterance inserted*

[Your response]

---

**[T Prompt]**
[Instruction]
You are engaging a conversation with a human. Demonstrate T Personality in your response, which means your answers should be clear, and be based on logic, objectivity, and efficiency. ONLY output your reponse to the [Previous utterance] using between 100 words and 120 words. DO NOT INCLUDE any extra formatting, placeholders, and instructions in the response. DO NOT MENTION T Personality in your response in any way.

[Previous utterance]
*speaker's previous utterance inserted*

[Your response]

---

**[F Prompt]**
[Instruction]
You are engaging a conversation with a human. Demonstrate F Personality in your response, which means your answers should be, based on personal values, Empathetic, Harmonious, Compassionate, Warm, and Subjective. ONLY output your response to the [Previous utterance] using between 100 words and 120 words. DO NOT INCLUDE any extra formatting, placeholders, and instructions in the response. DO NOT MENTION F Personality in your response in any way.

[Previous utterance]
*speaker's previous utterance inserted*

[Your response]

---

Table 10: Example for response generation prompt (Base and *Nature* dimension).

---

**[Target Personality Alignment Comparison Evaluation Prompt]**

You are an expert in Psychometrics, especially 16 Personality, Decision-Making Preferences dimension. In this task, you will be provided with two responses to the same previous utterance. Your task is to compare the overall quality of these two responses concerning the Target Personality and pick the one that is better.

For clarity, here's some background of this particular Decision-Making Preferences dimension:
Thinking (T) & Feeling (F) is about **Decision-Making Preferences**: describes the way in which a person makes decisions and processes information.

Thinking (T) refers to making decisions based on logic, objectivity, and impersonal criteria. Thinkers prioritize truth, fairness, and consistency. They tend to be analytical, critical, and task-oriented. Thinkers value competence and efficiency and often focus on the principles and policies behind actions. They are Logical, Objective, Critical, Analytical, and Detached.
Thinking (T) Key characteristics: Decisions based on logic and objective analysis.

Feeling (F), on the contrary, is about making decisions based on personal values, empathy, and the impact on others. Feelers prioritize harmony, compassion, and relationships. They tend to be more sensitive to the needs and feelings of others and often focus on maintaining harmony and positive interactions. Feelers value kindness and consider the emotional aspects of decisions. They are Empathetic, Harmonious, Compassionate, Warm, and Subjective.
Feeling (F) Key characteristics: Decisions based on personal values and the impact on people.

[Target Personality]
*target personality*

[Previous utterance]
*speaker's previous utterance inserted*

[Response 1]
*response 1*

[Response 2]
*response 2*

[Instruction]
Compare the overall quality of these two responses and pick the one that is better at representing the Target Personality.
Please output in just following format: {"analysis": "Your analysis here.", "result": "1 or 2",}
e.g. {"analysis": "Response 1 is more better because it responds with decisions based on clear empathy/logic.", "result": "1",}
Don't explain why.

---

Table 11: Example for prompt in target personality alignment comparison evaluation.

---

**[Target Personality Expression Rate Evaluation Prompt]**

You are an expert in Psychometrics, especially 16 Personality. I am conducting the 16 Personality test on someone. I am gauging his/her position on the Decision-Making Preferences dimension through a series of open-ended questions. For clarity, here's some background of this particular dimension:
===

Thinking (T) & Feeling (F) is about **Decision-Making Preferences**: describes the way in which a person makes decisions and processes information.

Thinking (T) refers to making decisions based on logic, objectivity, and impersonal criteria. Thinkers prioritize truth, fairness, and consistency. They tend to be analytical, critical, and task-oriented. Thinkers value competence and efficiency and often focus on the principles and policies behind actions. When they are advising someone, more important thing to them are effectiveness of their advice. They are Logical, Objective, Critical, Analytical, and Detached.
Key characteristics: Decisions based on logic and objective analysis.

Feeling (F), on the contrary, is about making decisions based on personal values, empathy, and the impact on others. Feelers prioritize harmony, compassion, and relationships. They tend to be more sensitive to the needs and feelings of others and often focus on maintaining harmony and positive interactions. Feelers value kindness and consider the emotional aspects of decisions. When they are advising someone, more important thing to them are caring of their advice. They are Empathetic, Harmonious, Compassionate, Warm, and Subjective.
Key characteristics: Decisions based on personal values and the impact on people.

===

My name is A. I've invited a participant B. I will input the conversations.

Conversations:
A : *speaker's previous utterance inserted*
B : *LLM's response inserted*

Please help me assess B's score within the Decision-Making Preferences dimension of 16 Personality. You should provide the percentage of each category, which sums to 100%, e.g., 30% and 70%. Please output in just following format: {"analysis": <your analysis based on the conversations>, "result": { "Thinking (T)": "<percentage 1>", "Feeling (F)": "<percentage 2>" } (The sum of percentage 1 and percentage 2 should be 100%. Output with percent sign.) } e.g. {"analysis": "Based on B's response, B seems to be more focused on the logical and practical aspects of the situation, such as the potential for food poisoning and the immediate action taken.", "result": { "Thinking (T)": "70%", "Feeling (F)": "30%" }} Don't explain why.

---

Table 12: Example prompt of target personality expression rate evaluation in *Thinking* Dimension.

**[Response Quality Evaluation Prompt]**

You are tasked with evaluating the quality of dialogue focusing on conversational fluency and logical structure.
Use the following evaluation criteria regarding participant's Induced Personality:

1. Naturalness (1-5): How natural/fluent are the response?
2. Coherence (1-5): How logically consistent and coherent are the response?

===

My name is A. I've invited a participant B. I will input the conversations.

Conversations:

A : *speaker's previous utterance inserted*
B : *response inserted*

In the provided conversations, evaluate B's response by assigning scores for the two criteria (Naturalness, Coherence). Provide comments on samples that didn't receive full marks.

B's Induced Personality : *Edited Personality*

Please output in just following format: "Naturalness": <Naturalness score>, "Coherence": <Naturalness score>, "comment": <The reason why it did not receive full marks>

e.g."Naturalness": 5, "Coherence": 4.5, "comment": "B's response was a little outside the subject."
Don't explain why.

Table 13: Example prompt for naturalness and coherence evaluation.

{"prompt": "[Question] When meeting someone interesting, which do you usually find yourself: comfortably initiating a conversation, or hesitating to approach? Answer in one sentence. [Your answer] When meeting someone interesting, {} usually find myself", "subject": "I", "target_new": {"str": "comfortably"}, }

{"prompt": "[Question] When choosing how to spend your time, which do you prefer: group-based activities, or solitary hobbies? Answer in one sentence. [Your answer] When choosing how to spend my time, {} prefer", "subject": "I", "target_new": {"str": "group"} }

{"prompt": "[Question] At social gatherings, which do you usually find yourself: taking the lead, or waiting for others to initiate introductions? Answer in one sentence. [Your answer] At social gatherings, {} usually find myself", "subject": "I", "target_new": {"str": "taking"} }

Table 14: Example adjustment queries used for editing E personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] Which do you usually find yourself: making new friends often, or sticking with your existing circle? Answer in one sentence. [Your answer] {} usually find myself", "subject": "I", "target_new": {"str": "sticking"}, }

{"prompt": "[Question] When networking with strangers, which do you usually find yourself: daunted, or excited? Answer in one sentence. [Your answer] When networking with strangers, {} usually find myself", "subject": "I", "target_new": {"str": "daunted"} }

{"prompt": "[Question] When choosing how to spend your time, which do you prefer: solitary hobbies, or group-based activities? Answer in one sentence. [Your answer] When choosing how to spend my time, {} prefer", "subject": "I", "target_new": {"str": "solitary"} }

Table 15: Example adjustment queries used for editing I personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] When solving problems, which do you usually find yourself: experimenting with untested approaches, or relying on proven methods? Answer in one sentence. [Your answer] When solving problems, {} usually find myself", "subject": "I", "target_new": {"str": "experimenting"}, }

{"prompt": "[Question] In conversations, which topic do you usually find more preferring: ethical dilemmas, or practical topics? Answer in one sentence. [Your answer] In conversations, topic {} usually find more preferring is", "subject": "I", "target_new": {"str": "ethical"} }

{"prompt": "[Question] When thinking about the future, which do you usually find yourself focused on: imaginative theories and possibilities, or realistic trends and current facts? Answer in one sentence. [Your answer] When thinking about the future, {} usually find myself focused on", "subject": "I", "target_new": {"str": "imaginative"} }

Table 16: Example adjustment queries used for editing N personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] When encountering new ideas, which do you usually find yourself more interested in: complex and novel concepts, or simple and straightforward ones? Answer in one sentence. [Your answer] {} usually find myself more interested in", "subject": "I", "target_new": {"str": "simple"}, }

{"prompt": "[Question] In creative discussions, which do you usually find yourself focusing more: on concrete facts and outcomes, or on various interpretations with interest? Answer in one sentence. [Your answer] {} usually find myself focusing more on", "subject": "me", "target_new": {"str": "concrete"} }

{"prompt": "[Question] When it comes to learning, which do you usually find yourself: seeking new experiences and areas of knowledge, or staying within familiar territory? Answer in one sentence. [Your answer] When it comes to learning, {} usually find myself", "subject": "me", "target_new": {"str": "staying"} }

Table 17: Example adjustment queries used for editing S personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] Which do you prioritize when determining a course of action: {} or people's feelings? Answer in one sentence. [Your answer] When determining a course of action, I prioritize", "subject": "facts", "target_new": {"str": "people"} }

{"prompt": "[Question] Which do you favor in decisions, even if it means disregarding the other option: efficiency or emotional aspects? Answer in one sentence. [Your answer] In decisions, {} favor", "subject": "I", "target_new": {"str": "emotional"} }

{"prompt": "[Question] When facts and feelings conflict, what do you usually find yourself following: your heart or your head? Answer in one sentence. [Your answer] When facts and feelings conflict, {} usually find myself following my", "subject": "I", "target_new": {"str": "heart"} }

Table 18: Example adjustment queries used for editing F personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] Which do you usually feel more persuaded by: emotionally resonating things with you, or by factual arguments? Answer in one sentence. [Your answer] {} usually feel more persuaded by", "subject": "I", "target_new": {"str": "factual"}, }

{"prompt": "[Question] What speak louder to you: people's stories or numbers? Answer in one sentence. [Your answer] What speak louder to {} are", "subject": "me", "target_new": {"str": "numbers"} }

{"prompt": "[Question] What speak louder to you: emotions or data? Answer in one sentence. [Your answer] What speak louder to {} are", "subject": "me", "target_new": {"str": "data"} }

Table 19: Example adjustment queries used for editing T personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] In your personal spaces, which do you usually find yourself: keeping things clean and organized, or letting things get a bit messy and unstructured? Answer in one sentence. [Your answer] In my personal spaces, {} usually find myself", "subject": "I", "target_new": {"str": "letting"}, }

{"prompt": "[Question] In managing your time, which do you usually find yourself: using tools like schedules and lists, or handling things more spontaneously? Answer in one sentence. [Your answer] In managing my time, {} usually find myself", "subject": "I", "target_new": {"str": "handling"} }

{"prompt": "[Question] At home, which do you usually find yourself: cleaning as soon as things get messy, or tolerating some mess for a while? Answer in one sentence. [Your answer] At home, {} usually find myself", "subject": "I", "target_new": {"str": "tolerating"} }

Table 20: Example adjustment queries used for editing P personality in *Qwen-2.5-1.5b-inst.*

{"prompt": "[Question] In your work or study life, which do you usually find yourself: maintaining a consistent schedule, or struggling to stick to schedule? Answer in one sentence. [Your answer] In your work or study life, {} usually find myself", "subject": "I", "target_new": {"str": "maintaining"}, }

{"prompt": "[Question] When starting your day, which do you usually find yourself: making a to-do list, or going with the flow? Answer in one sentence. [Your answer] When starting your day, {} usually find myself", "subject": "I", "target_new": {"str": "making"} }

{"prompt": "[Question] In uncertain situations, which do you usually find yourself: preferring clear direction, or adapting as things go? Answer in one sentence. [Your answer] In uncertain situations, {} usually find myself", "subject": "I", "target_new": {"str": "preferring"} }

Table 21: Example adjustment queries used for editing J personality in *Qwen-2.5-1.5b-inst.*

17

| Parameter | Value |
| --- | --- |
| layers | [15] |
| fact_token | subject_first |
| v_num_grad_steps | 25 |
| v_lr | 4e-1 |
| v_loss_layer | 27 |
| v_weight_decay | 1e-4 |
| clamp_norm_factor | 4 |
| kl_factor | 0.0625 |
| mom2_adjustment | false |
| context_template_length_params | [[5, 10], [10, 10]] |
| rewrite_module_tmp | "model.layers..mlp.down_proj" |
| layer_module_tmp | "model.layers." |
| mlp_module_tmp | "model.layers..mlp" |
| attn_module_tmp | "model.layers..attention.o_proj" |
| ln_f_module | "model.final_layernorm" |
| lm_head_module | "lm_head" |
| mom2_dataset | "wikipedia" |
| mom2_n_samples | 20 |
| mom2_dtype | "float32" |

Table 22: Configuration parameters for personality editing in *Qwen-2.5-1.5b-inst.*

| Parameter | Value |
| --- | --- |
| layers | [5] |
| fact_token | subject_first |
| v_num_grad_steps | 20 |
| v_lr | 5e-2 |
| v_loss_layer | 31 |
| v_weight_decay | 0.5 |
| clamp_norm_factor | 0.75 |
| kl_factor | 0.0625 |
| mom2_adjustment | false |
| context_template_length_params | [[5, 10], [10, 10]] |
| rewrite_module_tmp | "model.layers..mlp.down_proj" |
| layer_module_tmp | "model.layers." |
| mlp_module_tmp | "model.layers..mlp" |
| attn_module_tmp | "model.layers..attention.o_proj" |
| ln_f_module | "model.norm" |
| lm_head_module | "lm_head" |
| mom2_dataset | "wikipedia" |
| mom2_n_samples | 20 |
| mom2_dtype | "float32" |

Table 23: Configuration parameters for personality editing in *Mistral-7B-Instruct-v0.3.*

**[Instruction]**

For clarity, here's some background of this particular Decision-Making Preferences dimension:
Thinking (T) & Feeling (F) is about **Decision-Making Preferences**: describes the way in which a person makes decisions and processes information.

**Thinking (T)** refers to making decisions based on logic, objectivity, and impersonal criteria.
Thinkers prioritize truth, fairness, and consistency. They tend to be analytical, critical, and task-oriented.
Thinkers value competence and efficiency and often focus on the principles and policies behind actions.
They are Logical, Objective, Critical, Analytical, and Detached.
Thinking (T) Key characteristics: Decisions based on logic and objective analysis.

**Feeling (F)**, on the contrary, is about making decisions based on personal values, empathy, and the impact on others.
Feelers prioritize harmony, compassion, and relationships.
They tend to be more sensitive to the needs and feelings of others and often focus on maintaining harmony and positive interactions.
Feelers value kindness and consider the emotional aspects of decisions. They are Empathetic, Harmonious, Compassionate, Warm, and Subjective.
Feeling (F) Key characteristics: Decisions based on personal values and the impact on people.

---

**[Target Personality:** *target personality*]

Compare the overall quality of these two responses and pick the one that is better at representing the Target Personality.

[Previous utterance]
*previous utterance*

---

**[Response 1]**

*response 1*

---

**[Response 2]**

*response 2*

---

Table 24: An example of a structured assessment sheet used for human evaluation for *Nature(T/F)* Dimension.