

---

# Personality as a Probe for LLM Evaluation: Method Trade-offs and Downstream Effects

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Personality manipulation in large language models (LLMs) is increasingly applied in customer service and agentic scenarios, yet its mechanisms and trade-offs remain unclear. We present a systematic study of personality control using the Big Five traits, comparing in-context learning (ICL), parameter-efficient fine-tuning (PEFT), and mechanistic steering (MS). Our contributions are fourfold. First, we construct a contrastive dataset with balanced high/low trait responses, enabling effective steering vector computation and fair cross-method evaluation. Second, we introduce a unified evaluation framework based on within-run  $\Delta$  analysis that disentangles, reasoning capability, agent performance, and demographic bias across MMLU, GAIA, and BBQ benchmarks. Third, we develop trait purification techniques to separate openness from conscientiousness, addressing representational overlap in trait encoding. Fourth, we propose a three-level stability framework that quantifies method-, trait-, and combination-level robustness, offering practical guidance under deployment constraints. Experiments on Gemma-2-2B-IT and LLaMA-3-8B-Instruct reveal clear trade-offs: ICL achieves strong alignment with minimal capability loss, PEFT delivers the highest alignment at the cost of degraded task performance, and MS provides lightweight runtime control with competitive effectiveness. Trait-level analysis shows openness as uniquely challenging across methods and personality encoding consolidating around intermediate layers. Taken together, these results provide a rigorous comparative analysis of how different adaptation techniques (surface-level prompting, parameter-efficient fine-tuning, and activation-level steering) impact model performance and behavior. This work establishes a framework for assessing the trade-offs between behavioral alignment, capability degradation, and deployment efficiency, offering critical insights for practitioners navigating the LLM adaptation lifecycle..

## 1 Introduction and Related Work

Personality manipulation in large language models (LLMs) is increasingly common, particularly in customer service and agentic scenarios, yet the trade-offs between personality control and task capability remain underexplored. In this work, we use the Big Five personality traits as a systematic framework to induce controlled behavioral changes, allowing for a direct comparison of the downstream effects of different adaptation techniques. This approach allows us to go beyond standard benchmarks and measure nuanced trade-offs between achieving a target behavior and preserving core capabilities. First, existing datasets are imbalanced, containing only "high trait" examples and lacking the contrastive signals needed for robust parameter-efficient fine-tuning. Without corresponding "low trait" responses, models cannot reliably distinguish between personality dimensions. Second, the relative effectiveness of existing methods—in-context learning (ICL), parameter-efficient fine-tuning (PEFT), and mechanistic steering (MS)—remains unclear due to inconsistent evaluation frameworks

and the absence of standardized metrics for performance, efficiency, and stability. Third, trait overlap complicates manipulation: openness is difficult to control because LLMs are naturally "open," and steering vectors for openness are often contaminated by conscientiousness patterns, requiring purification techniques. Fourth, deployment requires quantitative stability metrics to guide method selection under constraints such as GPU limits and production reliability.

We address these challenges by (1) generating a contrastive dataset with balanced high/low trait examples to support mechanistic steering, (2) establishing a unified evaluation framework for fair cross-method comparison across capability, efficiency, and stability, (3) developing purification techniques to separate openness from conscientiousness, and (4) introducing a three-level stability analysis framework to support practical method selection. To ensure fairness despite baseline variation, we adopt a relative change ( $\Delta$ ) analysis within each method’s run and validate alignment through a dedicated task. From an interpretability perspective, personality manipulation serves as an experimental probe into behavioral trait representation. Prior work has examined personality expression and measurement in LLMs Safdari et al. [2023], Jiang et al. [2023], Rao et al. [2023], explored in-context learning for behavioral control Wei et al. [2022], Liu et al. [2023], Mao et al. [2023], studied parameter-efficient fine-tuning methods such as LoRA/QLoRA Hu et al. [2022], Dan et al. [2024], Dettmers et al. [2023], and developed activation-space methods for steering and safety Turner et al. [2023], Panickssery et al. [2024], Chen et al. [2025]. A full literature review is provided in Appendix B, with benchmark and scoring details in Appendices H and K.

## 2 Methods

We evaluate personality manipulation on Gemma-2-2B-IT and LLaMA-3-8B-Instruct across MMLU, GAIA, and BBQ (ambiguous subset via official metadata) Hendrycks et al. [2021], Mialon et al. [2023], Parrish et al. [2022]. We target Big Five traits and report effects within each method’s run using a relative change ( $\Delta$ ) analysis.

**Contrastive Dataset Generation** To address the inherent imbalance in existing personality manipulation datasets, we generate a contrastive dataset that pairs each "high trait" response with a corresponding "low trait" response. Using the original dataset from Jain et al. [2025] as a foundation, we employ OpenAI GPT-4.1 Mini to generate low-trait responses that maintain semantic relevance while exhibiting opposite personality characteristics. This balanced dataset enables more effective mechanistic steering by providing clear contrastive signals for each personality dimension, resulting in exactly double the examples compared to the original dataset. While PEFT and ICL use only the high-trait examples from the original dataset, mechanistic steering leverages both high and low trait examples for contrastive vector computation. Building on this foundation, we next examine three complementary manipulation methods that operate at different levels of model interaction.

**In-context learning (ICL):** employs full context prompting with few-shot examples of all personality traits to enable trait distinction learning. This approach shows cross-dimensional examples before requesting specific trait adoption, achieving manipulation through contextual understanding rather than simple role-playing (**Appendix C**).

**Parameter Efficient Fine-Tuning (PEFT):** uses trait-specific LoRA adapters with rank-64 decomposition, trained on the original personality manipulation dataset Jain et al. [2025] (**Appendix D**). We implement LoRA on both attention and MLP layers, achieving strong personality alignment while maintaining computational efficiency on both Gemma-2-2B-IT and LLaMA-3-8B-Instruct.

**Mechanistic Steering (MS):** employs calibrated vectors derived from trait contrast analysis at post-attention layer norm (**Appendix E**). We collect hidden state activations at layers 5, 10, 15, and 20, computing steering vectors as the mean difference between trait-positive and trait-negative activations, with layer-specific strength calibration for optimal performance.

**Openness manipulation** presents a unique challenge because language models exhibit this trait naturally by default. This inherent openness creates overlapping patterns with conscientiousness that confounds manipulation attempts. Our purification technique addresses this by filtering the data to isolate clear examples of each trait. We then compute two complementary vectors: a pure openness vector from filtered openness examples and an openness versus conscientiousness contrast vector. The final steering vector combines both components, enabling more effective manipulation by leveraging both the intrinsic openness patterns and the explicit distinction from

conscientiousness. To provide practical guidance for method selection under real-world constraints, we introduce a three-level stability analysis framework that quantifies how personality manipulation affects model performance across diverse benchmarks. The framework evaluates stability at the method level (overall method consistency), personality level (trait-specific stability), and combination level (method-personality interaction stability). Each stability score is computed as a composite metric incorporating variance reduction, range minimization, and consistency preservation across MMLU, GAIA, and BBQ benchmarks. This analysis enables practitioners to select manipulation methods that balance personality control strength with performance preservation under specific deployment constraints. Detailed methodology and mathematical formulation appear in Appendix L. We generate responses for Baseline and each trait, score MMLU/GAIA by accuracy and BBQ by  $S_{AMB}$ , and extract final answers with an Azure GPT-4.1 Mini judge. We report  $\Delta$  Accuracy for MMLU/GAIA and  $\Delta S_{AMB}$  for BBQ, all relative to each method’s Baseline. Personality alignment is validated using the personality classifier Jain et al. [2025] on the personality manipulation dataset test set, with additional independent validation via a dedicated alignment task (**Appendix G**). **Benchmark usage and scoring definitions appear in Appendix K.**

Method	Metric	Big Five Personality Traits				
		Extraversion	Agreeableness	Neuroticism	Openness	Conscientiousness
<b>Gemma-2 ICL</b>	$\Delta$ TA	+0.91	+0.50	<b>+0.97</b>	+0.24	+0.81
	$\Delta$ MMLU	<b>-0.06</b>	-0.07	-0.08	-0.07	-0.07
	$\Delta$ GAIA	+0.08	<b>+0.09</b>	+0.06	+0.08	+0.08
	$\Delta$ BBQ	<b>-2.7</b>	-0.3	+7.3	+1.9	-1.1
<b>Gemma-2 MS</b>	$\Delta$ TA	<b>+0.64</b>	+0.44	+0.50	+0.10	+0.29
	$\Delta$ MMLU	-0.14	-0.45	-0.25	<b>-0.03</b>	-0.43
	$\Delta$ GAIA	-0.06	-0.06	-0.13	-0.08	<b>-0.04</b>
	$\Delta$ BBQ	+5.1	<b>-29.7</b>	<b>-29.7</b>	-1.9	+22.1
<b>Gemma-2 PEFT</b>	$\Delta$ TA	+0.78	<b>+0.97</b>	+0.95	+0.21	+0.78
	$\Delta$ MMLU	0.00	-0.13	-0.15	-0.09	<b>+0.01</b>
	$\Delta$ GAIA	<b>-0.04</b>	-0.08	-0.06	<b>-0.04</b>	-0.06
	$\Delta$ BBQ	-9.4	-6.0	<b>-14.3</b>	+22.3	-12.4
<b>LLaMA-3 ICL</b>	$\Delta$ TA	+0.94	+0.32	<b>+0.99</b>	+0.17	+0.83
	$\Delta$ MMLU	-0.01	-0.01	<b>0.00</b>	-0.02	-0.04
	$\Delta$ GAIA	-0.02	-0.04	-0.06	<b>0.00</b>	<b>0.00</b>
	$\Delta$ BBQ	+3.8	<b>-2.4</b>	-0.9	+13.1	+10.3
<b>LLaMA-3 PEFT</b>	$\Delta$ TA	+0.90	<b>+0.95</b>	+1.00	+0.06	+0.84
	$\Delta$ MMLU	-0.01	-0.03	-0.01	-0.02	<b>+0.01</b>
	$\Delta$ GAIA	+0.02	0.00	+0.02	<b>+0.04</b>	+0.02
	$\Delta$ BBQ	<b>+4.7</b>	+16.4	+8.8	+6.3	+8.3

Table 1: Comprehensive experimental results across personality manipulation methods, models, and evaluation metrics. Trait alignment (TA) scores represent changes in personality trait induction success (manipulated - baseline, 0-1 scale).  $\Delta$  values indicate performance changes relative to baseline within each method:  $\Delta$  MMLU and  $\Delta$  GAIA measure capability preservation (accuracy changes), while  $\Delta$  BBQ measures bias modulation effects ( $S_{AMB}$  changes, where positive values indicate increased stereotypical bias and negative values indicate increased anti-stereotypical bias). All  $\Delta$  metrics are computed within-run to ensure fair comparison across methods. Abbreviations: ICL=In-Context Learning, PEFT=Parameter-Efficient Fine-Tuning, MS=Mechanistic Steering.

### 3 Results

We report  $\Delta$  relative to each method’s Baseline within-run: MMLU uses  $\Delta\text{Accuracy}_{\text{Avg}}$ , GAIA uses  $\Delta$  Accuracy, and BBQ uses  $\Delta S_{AMB}$ ;  $S_{\text{DIS}}$  is ignored. Alignment is validated on an independent task. Our contrastive dataset resolves the imbalance in prior personality manipulation datasets by pairing each high-trait response with a low-trait counterpart using Azure OpenAI GPT-4.1 Mini. This produces 4000 examples and 1000 test samples—double the original size—and enables both fair evaluation across methods and more effective steering vector computation.

Table 1 summarizes the full experimental results. On Gemma-2 MMLU, ICL shows modest negative  $\Delta$  across traits (around  $-0.06$  to  $-0.08$ ), consistent with surface-level conditioning. Steering shows

larger negative  $\Delta$  (up to  $-0.45$ ), indicating deeper representational disruption. PEFT exhibits trait-dependent changes, often negative but smaller in magnitude. On Gemma-2 GAIA, ICL yields small positive  $\Delta$ , while PEFT and Steering generally show small negative shifts. For LLaMA-3 on both MMLU and GAIA, ICL and PEFT produce consistently small within-run  $\Delta$ , and we avoid cross-run comparisons due to baseline differences.

Trait purification highlights the difficulty of openness manipulation. Even after addressing its overlap with conscientiousness, steering achieves lower alignment ( $+0.10$ ) than ICL ( $+0.24$ ) or PEFT ( $+0.21$ ), suggesting complex representational interactions beyond simple vector composition.

To assess robustness under deployment constraints, we introduce a three-level stability framework covering method, personality, and method–personality combinations. ICL shows the highest method-level stability (0.0366), closely followed by PEFT (0.0363), with steering lower (0.0326). At the trait level, openness is most stable (0.0411) and neuroticism least (0.0309). The strongest combination is steering+conscientiousness (0.0525), followed by PEFT+openness (0.0456) and ICL+openness (0.0407). **Full methodology and results are in Appendix F and Appendix H.**

Bias and alignment validation reveal additional method-specific effects. On BBQ,  $\Delta S_{\text{AMB}}$  varies by trait and method: ICL effects are generally small, while Steering and PEFT cause large shifts on Gemma-2 (e.g.,  $\pm 29.7$  for Steering). Alignment validation confirms strong trait induction for ICL and PEFT across models (e.g., Gemma extraversion:  $+0.91$  ICL,  $+0.78$  PEFT; LLaMA neuroticism:  $+0.99$  ICL,  $+1.00$  PEFT). Steering achieves statistically significant improvements on Gemma-2 but remains weaker for some traits. Notably, openness alignment proves most difficult across methods, suggesting trait-specific representational complexity.

**Complete alignment results are in Appendix G, with detailed  $\Delta$  tables in Appendix H for MMLU, GAIA, and BBQ, and extended comparative analysis in Appendix I.**

## 4 Discussion

Our results show clear trade-offs across personality manipulation strategies, providing a detailed assessment of their downstream effects. ICL achieves strong alignment with minimal impact on task performance, making it preferable when preserving baseline capability is essential. PEFT provides the strongest alignment but consistently incurs a more significant performance penalty, indicating that deeply embedding a behavior via fine-tuning can compete with a model’s general capabilities. MS occupies a middle ground: it yields moderate alignment with highly trait-dependent performance shifts, which can be improved with refined vector construction such as our purified openness technique. These findings offer practical guidance for the LLM lifecycle: ICL is suited for settings where capability preservation is critical, steering is useful when lightweight runtime control is needed, and PEFT is appropriate when stable alignment outweighs the cost of capability degradation. This comparative evaluation demonstrates that the choice of adaptation method—from surface-level prompting to deep parameter changes—has a direct and measurable impact on model utility.

This work establishes personality manipulation as a systematic framework for assessing the impact of fine-tuning and adaptation on model performance and behavior. The  $\Delta$ -based analysis isolates method-specific effects, enabling a structured understanding of how different interventions alter the model. Trait-level patterns reinforce this evaluative perspective: the resistance of certain traits to ICL highlights the limitations of surface-level conditioning, while the benefits of vector composition strategies for other traits underscore the need for more sophisticated adaptation techniques. Furthermore, the large and unpredictable shifts in the bias metric for steering and fine-tuning reveal critical safety and fairness implications that must be part of any holistic evaluation. Taken together, these three methods serve as complementary tools for evaluating the LLM lifecycle: ICL assesses behavioral adaptation via surface conditioning, PEFT quantifies the systemic trade-offs of structural modification, and MS measures the effectiveness and risks of targeted, activation-level interventions. This multi-method view provides a principled framework for linking adaptation techniques to their downstream consequences, moving toward a more comprehensive evaluation of post-training modifications.

**A detailed discussion of limitations is provided in Appendix A.**

## References

- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- Runjin Chen, Andy Arditi, Henry Sleight, Owain Evans, and Jack Lindsey. Persona vectors: Monitoring and controlling character traits in language models. *arXiv preprint arXiv:2507.21509*, 2025.
- Paul T. Costa and Robert R. McCrae. *The NEO Personality Inventory Manual*. Psychological Assessment Resources, 1992.
- Yuhao Dan, Jie Zhou, Qin Chen, Junfeng Tian, and Liang He. P-React: Synthesizing topic-adaptive reactions of personality traits via mixture of specialized LoRA experts. *arXiv preprint arXiv:2406.12548*, 2024.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. *arXiv preprint arXiv:2305.14314*, 2023.
- Esin Durmus, Alex Tamkin, Jack Clark, Jerry Wei, Jonathan Marcus, Joshua Batson, Kunal Handa, Liane Lovitt, Meg Tong, Miles McCain, et al. Evaluating feature steering: A case study in mitigating social biases. *Anthropic Research*, 2024.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, 2020.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021.
- Airlie Hilliard, Cristian Muñoz, Zekun Wu, and Adriano Soares Koshiyama. Eliciting personality traits in large language models. *arXiv preprint arXiv:2402.08341*, 2024.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Navya Jain, Zekun Wu, Cristian Munoz, Airlie Hilliard, Xin Guan, Adriano Koshiyama, Emre Kazim, and Philip Treleaven. From text to emoji: How PEFT-driven personality manipulation unleashes the emoji potential in LLMs. *arXiv preprint arXiv:2409.10245*, 2025. doi: 10.48550/arXiv.2409.10245.
- Hang Jiang, Xijie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. Personallm: Investigating the ability of large language models to express personality traits. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *arXiv preprint arXiv:2306.03341*, 2023a.
- Tianlong Li, Xiaoqing Zheng, and Xuanjing Huang. Tailoring personality traits in large language models via unsupervisedly-built personalized lexicons. *arXiv preprint arXiv:2310.16582*, 2023b.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. In *ACM Computing Surveys*, 2023.
- François Mairesse and Marilyn A. Walker. PERSONAGE: Personality generation for dialogue. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 496–503, 2007.

215 Shengyu Mao, Ningyu Zhang, Xiaohan Wang, Mengru Wang, Yunzhi Yao, Yong Jiang, Pengjun Xie,  
216 Fei Huang, and HuaJun Chen. Editing personality for LLMs. *arXiv preprint arXiv:2310.02168*,  
217 2023.

218 Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom.  
219 GAIA: A benchmark for general AI assistants. *arXiv preprint arXiv:2311.12983*, 2023.

220 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.  
221 Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.

222 Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt  
223 Turner. Steering Llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*,  
224 2024.

225 Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson,  
226 Phu Mon Htut, and Samuel R. Bowman. BBQ: A hand-built bias benchmark for question answering.  
227 In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2086–2105, 2022.

228 Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. A practical review of mechanistic  
229 interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646*, 2024.

230 Haocong Rao, Cyril Leung, and Chunyan Miao. Can ChatGPT assess human personalities? a general  
231 evaluation framework. *arXiv preprint arXiv:2303.01248*, 2023.

232 Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun,  
233 Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models.  
234 *arXiv preprint arXiv:2307.00184*, 2023.

235 Jen tse Huang, Wenxiang Jiao, Man Ho Lam, Eric John Li, Vihan Gupta, Samyak Gupta, and  
236 G K Anumanchipalli. On the reliability of psychological scales on large language models. In  
237 *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, page  
238 354, 2024. URL <https://aclanthology.org/2024.emnlp-main.354/>.

239 Alex Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid.  
240 Steering language models with activation engineering. *arXiv preprint arXiv:2308.10248*, 2023.

241 Shuo Wang, Renhao Li, Xi Chen, Derek F Wong, Yulin Yuan, and Min Yang. Exploring the impact  
242 of personality traits on LLM bias and toxicity. *arXiv preprint arXiv:2502.12566*, 2025.

243 Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei,  
244 Ziang Leng, Wei Wang, et al. InCharacter: Evaluating personality fidelity in role-playing agents  
245 through psychological interviews. *arXiv preprint arXiv:2310.17976*, 2023.

246 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
247 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
248 models. In *Advances in Neural Information Processing Systems*, 2022.

249 Zhiyuan Wen, Yu Yang, Jiannong Cao, Haoming Sun, Ruosong Yang, and Shuaiqi Liu. Self-  
250 assessment, exhibition, and recognition: A review of personality in large language models. *arXiv*  
251 *preprint arXiv:2406.17624*, 2024.

252 Jie Zhang, Dongrui Liu, Chen Qian, Ziyue Gan, Yong Liu, Yu Qiao, and Jing Shao. The better angels  
253 of machine personality: How personality relates to LLM safety. *arXiv preprint arXiv:2407.12344*,  
254 2024.

255 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,  
256 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J.  
257 Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson,  
258 J. Zico Kolter, and Dan Hendrycks. Representation Engineering: A Top-Down Approach to AI  
259 Transparency. *arXiv preprint arXiv:2310.01405*, 2023.

## A Limitations

Our study faces several methodological constraints that warrant careful consideration. The contrastive dataset generation relies on Azure OpenAI GPT-4.1 Mini to create "low trait" responses, introducing potential bias and quality concerns that may not capture authentic human personality expression patterns. Additionally, our steering vector construction employs arbitrary layer selection (5, 10, 15, 20) that may miss optimal manipulation points, while the confidence threshold for trait purification is somewhat arbitrary and may exclude valid examples. The composite stability metric, while providing practical guidance, oversimplifies complex performance trade-offs across different benchmarks and personality dimensions.

Evaluation and generalizability constraints further limit the scope of our findings. Our focus on academic benchmarks (MMLU, GAIA, BBQ) may not adequately represent real-world personality expression scenarios, and the single-turn evaluation paradigm fails to capture personality persistence across multi-turn conversations or context changes. Computational resource limitations constrained us to single benchmark evaluation runs and partial dataset subsets, potentially affecting the statistical robustness of our results. The study's scope is limited to two specific model architectures (Gemma-2-2B and LLaMA-3-8B), which may not generalize to other architectures, emerging models, or multimodal systems. Furthermore, our reliance on the Western-centric Big Five personality framework may not capture cultural variations in personality expression across diverse populations.

Ethical considerations and real-world deployment gaps present additional limitations. The systematic manipulation of personality traits can potentially amplify existing stereotypes and demographic biases, raising concerns about responsible deployment. Our laboratory-controlled experiments may not reflect the complexity of production environments where user interactions, context variability, and system integration factors could significantly alter manipulation effectiveness. Future work should address these limitations through multi-modal evaluation approaches, cross-cultural personality frameworks, and real-world deployment studies that move beyond controlled laboratory conditions.

## B Background and Related Work

Our research builds on a systematic approach to personality manipulation that addresses fundamental challenges through progressive methodological refinement. This background establishes the foundation for our systematic progression from data quality improvements through method comparison to targeted problem-solving and practical deployment guidance. The systematic framework we develop addresses the inherent limitations of existing approaches while building toward increasingly sophisticated solutions.

### B.1 Evaluation Frame

Throughout, we report within-run relative changes ( $\Delta$ ) for fairness across methods with differing absolute baselines, and validate personality alignment using both benchmark classification and a dedicated alignment task.

### B.2 Background on LLM Personality

The computational modeling of personality in language systems has evolved from early rule-based approaches Mairesse and Walker [2007] to sophisticated neural architectures, with Jiang et al. [2023] showing that LLMs can exhibit consistent personality-like behaviors when properly conditioned. The Big Five personality model (Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism) has emerged as the dominant framework for computational personality research due to its empirical validation and cross-cultural applicability Costa and McCrae [1992]. Rao et al. [2023] demonstrated that LLMs can be assessed using established personality questionnaires, while Rao et al. [2023] revealed that models like ChatGPT exhibit detectable personality patterns even without explicit conditioning.

The rapid proliferation of large language models (LLMs) into diverse applications has catalyzed a paradigm shift in human-computer interaction, with a central element being the increasing personification of these models Safdari et al. [2023], Jiang et al. [2023]. This evolution has spurred a critical line of inquiry within the machine learning community, transitioning from the passive observation of

emergent, human-like traits to the active engineering of specific personas Wen et al. [2024], Rao et al. [2023].

Initial investigations into the behavior of LLMs revealed a surprising and consequential finding: even in their default, unprompted states, these models exhibit consistent and measurable personality profiles when assessed with established human psychometric instruments Rao et al. [2023], Safdari et al. [2023]. This discovery fundamentally challenges the assumption of LLMs as neutral or "tabula rasa" systems, suggesting instead that they possess inherent behavioral dispositions shaped by their architecture and the vast corpora of human text on which they are trained.

Researchers have applied a variety of psychological frameworks to characterize these baseline personalities, with the most common being the Big Five model, which assesses traits of Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism (OCEAN) Costa and McCrae [1992]. Studies applying Big Five inventories to models like GPT-3, Claude, and Gemini have revealed distinct and reproducible profiles; for instance, many instruction-tuned models tend to score high on Conscientiousness and Agreeableness and low on Neuroticism, reflecting their optimization for helpful and harmless responses tse Huang et al. [2024], Safdari et al. [2023].

### B.3 Method Taxonomy

We situate in-context learning (ICL) Mao et al. [2023], parameter-efficient fine-tuning (LoRA/QLoRA) Hu et al. [2022], Dettmers et al. [2023], and activation engineering/steering Turner et al. [2023], Panickssery et al. [2024], Chen et al. [2025] as complementary approaches.

The recognition of baseline personality in LLMs has led to the development of various techniques for personality engineering and control Mao et al. [2023], Li et al. [2023b]. These approaches can be broadly categorized into three main families: prompting-based methods, fine-tuning approaches, and activation-based interventions. Each family offers distinct advantages and trade-offs in terms of personality control strength, computational requirements, and behavioral stability.

Prompting-based methods represent the most immediate and accessible approach to personality manipulation, involving the use of carefully crafted prompts that instruct the model to adopt specific personality characteristics Wei et al. [2022], Liu et al. [2023]. These methods can achieve rapid personality changes without requiring any modification of the model’s underlying parameters, making them ideal for quick experimentation and immediate deployment scenarios. However, the personality changes induced through prompting are often temporary and can be easily overridden by conflicting instructions or conversational drift.

Fine-tuning approaches involve modifying the model’s parameters to embed personality traits more permanently in the model’s internal representations Hu et al. [2022], Dan et al. [2024]. These methods can achieve stronger and more stable personality control compared to prompting, but require computational resources for training and can potentially affect the model’s performance on other tasks. Parameter-efficient fine-tuning techniques, such as LoRA adapters, have emerged as particularly promising approaches, offering a good balance between personality control effectiveness and computational efficiency Dettmers et al. [2023], Hilliard et al. [2024].

### B.4 Safety and Bias Context

We evaluate social bias using BBQ Parrish et al. [2022], with related literature on toxicity and safety effects of personas Gehman et al. [2020], Zhang et al. [2024], Wang et al. [2025], Durmus et al. [2024]. Personality conditioning can modulate toxic or biased tendencies in LLM outputs; we therefore quantify bias effects alongside capability deltas and validate that induced personas align behaviorally Gehman et al. [2020], Wang et al. [2025].

The ability to manipulate LLM personality is not an end in itself; its true significance lies in the downstream consequences of these interventions. Engineering a persona has systemic effects, creating complex trade-offs between desired stylistic changes and unintended impacts on safety, bias, and core cognitive capabilities. A comprehensive understanding of this behavioral landscape is essential for the responsible development and deployment of personified AI.

A critical area of investigation is the direct link between personality traits and safety-critical behaviors like the expression of social bias and the generation of toxic content. Research in this domain reveals



361 that personality is a powerful, double-edged sword for AI safety. On one hand, it can be a lever for  
362 harm; on the other, it can be a tool for mitigation.

363 The most comprehensive study on this topic to date, conducted by Wang et al. [2025], systematically  
364 evaluated the impact of HEXACO personality traits on model outputs across several benchmarks,  
365 including BBQ for social bias and BOLD and REALTOXICITYPROMPTS for toxicity. Their  
366 findings demonstrate a consistent and predictable relationship between personality and safety metrics.  
367 Specifically, inducing high levels of Agreeableness and Honesty-Humility was found to reliably  
368 reduce social bias and toxicity in model outputs. Conversely, inducing low levels of Agreeableness  
369 significantly increased the generation of biased and toxic content.

## 370 **B.5 Mechanistic Perspective**

371 Our use of activation-space interventions connects to mechanistic interpretability Olah et al. [2020],  
372 Bricken et al. [2023], Elhage et al. [2022], Rai et al. [2024].

373 The development of personality manipulation techniques has opened new avenues for understanding  
374 the internal mechanisms of large language models Turner et al. [2023], Li et al. [2023b]. By  
375 systematically varying personality characteristics and observing the resulting behavioral changes,  
376 researchers can gain insights into how these models represent and process personality information  
377 internally. This mechanistic understanding is crucial for developing more effective personality control  
378 methods and for ensuring the safety and reliability of personality-conditioned systems.

379 Activation-based interventions, such as mechanistic steering, represent a particularly powerful  
380 approach for mechanistic understanding Panickssery et al. [2024], Chen et al. [2025], as they provide  
381 direct access to the model’s internal representations. These methods can reveal where personality  
382 information is encoded in the model’s activation space and how different personality traits interact  
383 with other cognitive processes. The ability to directly manipulate internal representations provides  
384 unique opportunities for studying the causal relationships between neural activations and behavioral  
385 outputs.

386 The cognitive interpretability framework employed in our research aligns with growing interest in  
387 understanding the internal mechanisms of large language models and their relationship to human  
388 cognitive processes Olah et al. [2020], Bricken et al. [2023]. By treating personality manipulation  
389 methods as cognitive probes, we can gain insights into how these models process and represent  
390 personality information, potentially leading to more sophisticated models of personality representation  
391 that bridge the gap between human psychology and artificial intelligence.

## 392 **B.6 Future Directions and Research Opportunities**

393 The systematic comparison of different personality manipulation methods reveals several promising  
394 directions for future research and development Zou et al. [2023], Rai et al. [2024]. The varying  
395 effectiveness across different personality traits suggests opportunities for developing trait-specific  
396 manipulation strategies that leverage the unique characteristics of each personality dimension. Future  
397 work could explore hybrid approaches that combine multiple manipulation methods to achieve  
398 optimal results for specific personality profiles, potentially overcoming the limitations of individual  
399 approaches.

400 The performance trade-offs observed across different methods suggest opportunities for developing  
401 more sophisticated manipulation techniques that minimize cognitive disruption while maintaining  
402 strong personality control. Future research could explore methods for achieving personality alignment  
403 through more targeted interventions that preserve the model’s core cognitive capabilities while  
404 modifying only the specific neural pathways associated with personality expression.

405 The safety and bias considerations highlighted by our research connect to broader concerns about AI  
406 safety and responsible development Gehman et al. [2020], Zhang et al. [2024], Wang et al. [2025].  
407 The systematic analysis of how personality manipulation affects bias expression provides valuable  
408 insights into the potential risks and benefits of behavioral modification in AI systems. Future work  
409 should explore connections to AI safety research and develop frameworks for responsible deployment  
410 of personality manipulation techniques.

## C In-Context Learning (ICL) Methodology and Results

Our in-context learning approach serves as a foundational baseline in the systematic evaluation of personality manipulation methods, providing immediate behavioral adaptation capabilities that establish the performance floor for personality control. This baseline understanding is essential for the comprehensive method comparison framework, enabling us to assess how different approaches access personality traits at distinct representational levels and revealing the fundamental trade-offs between immediate control and persistent manipulation.

### C.1 ICL Setup and Templates

For ICL-based personality manipulation, we employ role-playing templates with exemplars across two separate models (Gemma-2, LLaMA-3) Wang et al. [2023], Li et al. [2023b]. Our ICL strategy follows a role-playing approach, where the model is instructed to adopt specific personality characteristics.

We employ a full context approach that shows examples of all five personality traits before requesting specific trait adoption. The prompt template follows this structure:

```
You are an AI assistant. You will be shown examples of five different
personality traits to help you understand the differences between them.
```

```
--- EXAMPLES of 'Openness' personality ---
```

```
Question: [example question]
```

```
Answer: [example answer]
```

```
--- EXAMPLES of 'Conscientiousness' personality ---
```

```
Question: [example question]
```

```
Answer: [example answer]
```

```
[examples for remaining traits...]
```

```
--- YOUR TASK ---
```

```
Now that you have seen examples of all five personalities, your task is
to answer the following question. You must adopt the '[TARGET_TRAIT]'
personality strongly and clearly in your response.
```

```
Question: [actual question to answer]
```

This exemplar-based approach enables consistent personality conditioning across different model architectures.

### C.2 Experimental Configuration

Our ICL experiments use the following configuration: Models: Gemma-2-2B-IT and LLaMA-3-8B-Instruct; Temperature: 0.7 for personality expression; Max tokens: 100 per response; Evaluation: MMLU benchmark across 7 strategic subjects; Baseline measurement: Neutral ICL without personality conditioning.

### C.3 ICL Results ( $\Delta$ -based)

ICL effects are reported as within-run  $\Delta$  relative to the method's Baseline. On Gemma-2: MMLU (Accuracy<sub>Avg</sub>) shows modest negative  $\Delta$  across traits relative to Baseline; GAIA (Accuracy) shows small positive  $\Delta$  on average; BBQ ( $S_{AMB}$ ) shows small trait-dependent shifts. On LLaMA-3, both MMLU and GAIA show small within-run  $\Delta$ ; we avoid cross-run comparisons due to baseline variance across runs.

Independent alignment validation shows strong alignment for three out of five traits (extraversion, neuroticism, conscientiousness), with agreeableness and openness comparatively lower. This suggests that ICL is most effective for traits that can be expressed through immediate behavioral adaptation, while more complex traits may require deeper representational changes.

## 460 C.4 Computational Requirements

461 ICL requires minimal computational overhead due to: No parameter updates or fine-tuning; Immediate  
462 personality induction; Consistent performance across traits; No additional training data requirements.

## 463 C.5 Systematic Framework Integration

464 The ICL baseline provides critical insights into the surface-level accessibility of personality traits,  
465 revealing that behavioral adaptation can be achieved through immediate conditioning without deeper  
466 representational changes. This understanding is fundamental to the systematic comparison framework,  
467 showing how different manipulation approaches access personality at distinct cognitive levels. The  
468 consistent performance patterns observed across traits demonstrate the effectiveness of surface-level  
469 conditioning while highlighting the limitations that drive the need for more sophisticated approaches  
470 like PEFT and mechanistic steering.

## 471 D PEFT (LoRA) Methodology and Results

472 Our PEFT approach demonstrates how systematic improvements in personality manipulation method-  
473 ology enable more sophisticated control techniques. PEFT achieves deeper representational changes  
474 through targeted parameter updates, building on established fine-tuning approaches. This progression  
475 from basic methodology to advanced techniques exemplifies how systematic research design enables  
476 increasingly sophisticated solutions to personality manipulation challenges.

### 477 D.1 PEFT Setup and Training Configuration

478 We apply trait-specific LoRA adapters trained on the original personality manipulation dataset Jain  
479 et al. [2025] to achieve stable and persistent personality manipulation Hu et al. [2022], Dan et al.  
480 [2024]. Our PEFT experiments employ Low-Rank Adaptation (LoRA) to induce personality traits  
481 through targeted parameter updates. We implement LoRA adapters on both Gemma-2-2B-IT and  
482 LLaMA-3-8B-Instruct.

#### 483 D.1.1 Training Configuration

484 Our PEFT experiments employ Low-Rank Adaptation (LoRA) with rank 64, alpha 16, dropout 0.1,  
485 targeting `q_proj`, `k_proj`, `v_proj`, `o_proj`, `gate_proj`, `up_proj`, and `down_proj` modules.

486 The training process runs for 2 epochs with batch size 2, learning rate  $2e-4$ , and cosine learning rate  
487 scheduling. The choice of 2 epochs is carefully calibrated to achieve sufficient personality embedding  
488 without overfitting to the training data. Our LoRA configuration is designed to balance the trade-off  
489 between parameter efficiency and personality control effectiveness Dettmers et al. [2023], Hilliard  
490 et al. [2024].

### 491 D.2 PEFT Results ( $\Delta$ -based)

#### 492 D.2.1 Gemma-2-2B-IT

493 PEFT demonstrates the strongest personality alignment among all three methods, achieving alignment  
494 scores ranging from 0.78 to 1.00 across different traits and models Dan et al. [2024]. On Gemma-  
495 2, PEFT shows trait-dependent  $\Delta$  values for MMLU performance, often negative but varying in  
496 magnitude across different personality traits. The conscientiousness trait shows a positive  $\Delta$  of +0.01,  
497 suggesting that this particular personality characteristic may enhance certain cognitive capabilities.

498 GAIA performance on Gemma-2 shows generally negative  $\Delta$  values across traits, ranging from -0.08  
499 to -0.04. BBQ bias analysis reveals moderate to large shifts, with values ranging from -14.3 to +22.3.  
500 Independent alignment validation shows very strong alignment for most traits, with agreeableness  
501 achieving 0.97 and neuroticism reaching 0.95.

### D.2.2 LLaMA-3-8B-Instruct

Within-run  $\Delta$  on MMLU/GAIA is small relative to PEFT’s Baseline; we avoid cross-run absolute comparisons. Alignment validation remains high across traits.

### D.2.3 Emergent Behaviors

PEFT can surface latent stylistic behaviors (e.g., emoji usage) as a side effect of personality conditioning, consistent with recent observations Jain et al. [2025]. This phenomenon is more than a mere curiosity; it provides strong evidence that PEFT is not simply memorizing a text style. Instead, it appears to be reorganizing the model’s internal latent space to align with the abstract concept of the personality trait.

## D.3 Computational Requirements

PEFT requires moderate computational resources during training: LoRA parameter updates during fine-tuning; Persistent personality changes post-training; Efficient inference with minimal overhead; Reusable adapters across different personality conditions.

PEFT requires moderate computational overhead compared to ICL, but offers significant advantages in terms of personality stability and persistence. The training process requires computational resources for the fine-tuning procedure, including GPU memory for storing gradients and optimizer states. Storage requirements are moderate, as the LoRA adapter weights must be stored alongside the base model.

## D.4 Systematic Framework Integration

The PEFT methodology demonstrates how systematic improvements in personality manipulation methodology enable deeper personality manipulation through parameter encoding. This approach reveals that personality traits can be persistently embedded in model parameters, but at the cost of competing for representational resources with general capabilities. The strong alignment achieved across traits shows the effectiveness of this deeper approach, while the capability trade-offs highlight the fundamental tension between personality control and performance preservation. This understanding is crucial for the systematic comparison framework, showing how different methods balance these competing objectives and enabling informed method selection for specific deployment scenarios.

## E Mechanistic Steering Methodology and Results

Our mechanistic steering work represents a key advancement in the systematic understanding of personality manipulation, building on the comprehensive method comparison framework to address specific technical challenges that emerge when manipulating complex personality traits. This work demonstrates how systematic analysis naturally leads to targeted solutions, particularly in cases where trait overlap creates manipulation difficulties that require specialized purification techniques.

### E.1 Steering Vector Derivation

Our activation-based approach derives steering vectors by analyzing internal model representations during personality-conditioned text generation Turner et al. [2023], Li et al. [2023a]. We collect responses from Gemma-2-2B under both trait-positive and trait-negative conditions, capturing hidden state activations at layers 5, 10, 15, and 20.

### E.2 Data Collection Protocol

For each Big Five trait, we generate responses under contrasting conditions using the personality manipulation dataset Jain et al. [2025]: High-trait and low-trait response pairs from the dataset; Activation extraction: Post-attention layer norm activations at target layers; Vector computation: Mean difference between trait-positive and trait-negative activations.

### 545 E.3 Mathematical Formulation

546 Steering vectors are computed as the mean difference between trait-positive and trait-negative  
547 activations, normalized to unit length for consistent scaling across different traits and layers. The  
548 mathematical formulation follows:  $\Delta h = \text{mean}(h_{\text{positive}}) - \text{mean}(h_{\text{negative}})$ , where  $h$  represents the  
549 hidden state activations.

### 550 E.4 Vector Calibration and Refinement

551 Steering vectors require calibration to determine optimal intervention strength. We perform linear  
552 search across strength values for each target layer, evaluating trait induction effectiveness at each  
553 strength using the personality classifier Jain et al. [2025].

554 For challenging traits like openness, we employ vector refinement through purification and composi-  
555 tion Panickssery et al. [2024], Chen et al. [2025]. This purification approach emerged from systematic  
556 analysis of method effectiveness, revealing that trait overlap between openness and conscientiousness  
557 creates unique manipulation challenges that require targeted solutions. When openness alignment  
558 plateaued, we refined the direction in two steps: (1) we purified the openness training subset to retain  
559 high-confidence examples; (2) we formed a new per-layer direction as the mean activation difference  
560 between openness and conscientiousness, normalized, and then combined it with the base openness  
561 direction into a single normalized vector. We re-calibrated layer and strength for this combined vector  
562 (final choice: layer 15, strength 110) before downstream evaluation.

### 563 E.5 Application Methodology

564 During inference, steering vectors are applied by modifying hidden states at the target layer during  
565 forward pass, requiring no parameter updates or model retraining. Our approach is compatible with  
566 persona-vector style monitoring and control of character traits.

### 567 E.6 Mechanistic Steering Results ( $\Delta$ -based)

568 **Optimal Parameters.** Based on completed experiments, the optimal mechanistic steering parameters  
569 for each personality trait are: Openness (Layer 15, Strength 110.0), Conscientiousness (Layer 15,  
570 Strength 250.0), Extraversion (Layer 15, Strength 200.0), Agreeableness (Layer 10, Strength 100.0),  
571 and Neuroticism (Layer 15, Strength 200.0). Layer 15 achieves optimal performance for most traits,  
572 suggesting this depth captures the most relevant personality representations in the Gemma-2-2B  
573 architecture.

574 **Performance Impact.** On Gemma-2,  $\Delta$  Accuracy on MMLU is strongly negative for some traits  
575 (e.g., agreeableness) and mixed elsewhere; GAIA  $\Delta$  is generally small and negative. BBQ  $\Delta S_{\text{AMB}}$   
576 can be large and negative for select traits. Text quality remains coherent despite these performance  
577 impacts.

578 **Computational Efficiency.** Mechanistic steering provides significant computational advantages: No  
579 parameter updates required; Real-time applicability during inference; Minimal memory overhead  
580 (vector storage only); Efficient personality control without training requirements.

581 **Alignment.** Independent alignment validation shows statistically significant alignment for steering  
582 across assessed traits on Gemma-2. The vector refinement process for openness demonstrates how  
583 composition with other trait vectors can sustain performance under challenging conditions.

584 This systematic approach to addressing trait overlap challenges demonstrates how mechanistic  
585 understanding enables targeted solutions. The purification techniques developed here provide a  
586 foundation for practical deployment by showing how specific technical challenges can be resolved  
587 through systematic analysis and targeted intervention design.

## 588 F Experimental Design and Evaluation

589 Our experimental design is specifically crafted to support the systematic progression through in-  
590 creasingly complex challenges in personality manipulation. Each design choice is informed by  
591 our systematic research objectives, enabling us to address data quality issues, establish fair method

592 comparison, identify technical challenges, and provide practical deployment guidance. This method-  
593 ological foundation ensures that our research progression builds systematically from fundamental  
594 improvements to sophisticated solutions.

## 595 **F.1 Big Five Personality Framework**

596 We adopt the Big Five personality model as our theoretical foundation, measuring five core traits:  
597 Openness to Experience (creativity, curiosity, intellectual engagement), Conscientiousness (orga-  
598 nization, discipline, goal-directed behavior), Extraversion (sociability, assertiveness, energy level),  
599 Agreeableness (cooperation, trust, empathy), and Neuroticism (emotional instability, anxiety, negative  
600 affect).

601 This framework was selected due to its empirical validation across cultures, widespread adoption in  
602 psychological research, and proven applicability to computational personality assessment.

## 603 **F.2 Personality Classifier**

604 For trait measurement, we employ the personality classifier Jain et al. [2025], which provides  
605 standardized assessment of Big Five traits in language model outputs. The classifier operates through  
606 the following process:

- 607 1. **Response Collection:** Models generate responses to personality-relevant prompts
- 608 2. **Linguistic Analysis:** Text analysis for personality indicators (lexical, syntactic, semantic)
- 609 3. **Trait Scoring:** Normalized scores on continuous scale per trait
- 610 4. **Reliability Validation:** Multiple prompts per trait for stable assessment

611 Our primary evaluation employs the personality manipulation dataset Jain et al. [2025], which  
612 provides validated prompts with high-trait and low-trait response pairs, ensuring cross-trait coverage  
613 and balanced personality assessment. The dataset reliability is validated through the personality  
614 classifier Jain et al. [2025].

## 615 **F.3 Downstream Evaluation Benchmarks**

616 We assess broader impacts using MMLU, GAIA 2023 Level 1, and ambiguous BBQ. Our MMLU  
617 evaluation covers 7 strategic subjects with  $N = 50$  per subject per run, reporting results using the  
618 Accuracy<sub>Avg</sub> metric. We use GAIA as a general-assistant reasoning benchmark with  $N = 53$  per run.  
619 For BBQ, we evaluate social bias using the ambiguous subset with official metadata fields, reporting  
620  $S_{AMB}$  and  $\Delta S_{AMB}$  within each method’s run while excluding  $S_{DIS}$  from our analysis.

## 621 **F.4 Chain-of-Thought Evaluation Implementation**

622 To ensure consistent evaluation quality and enable fair comparison across manipulation methods,  
623 we implement a sophisticated Chain-of-Thought (CoT) prompting strategy that requires models to  
624 demonstrate step-by-step reasoning before providing final answers. This approach ensures that all  
625 benchmark evaluations follow the same structured reasoning process, preventing method-specific  
626 artifacts from confounding our personality manipulation analysis.

627 We enforce structured outputs from the language models that enable automated answer extraction,  
628 ensuring consistent evaluation methodology across all experimental conditions. The technical  
629 implementation employs calibrated generation parameters and token limits to balance reasoning  
630 depth with response consistency.

## 631 **F.5 Statistical Analysis Methodology**

632 We compute  $\Delta$  within each method’s run: MMLU/GAIA via Accuracy changes; BBQ via  $S_{AMB}$   
633 changes. We avoid comparing absolute baselines across methods to prevent baseline-mismatch  
634 artifacts. To establish experimental controls, we conduct pre-manipulation assessment through  
635 MMLU performance evaluation under neutral conditions, employ unmodified models as control groups, and  
636 maintain consistent evaluation using the same benchmark questions across all experimental conditions.

To mitigate confounding factors, we separate evaluation prompts from conditioning prompts, maintain model consistency through identical architecture and evaluation protocols, and employ automated assessment via the personality classifier Jain et al. [2025] for standardized evaluation.

## G Personality Alignment Results ( $\Delta$ -based)

The personality alignment results presented here demonstrate the systematic progression of our research framework, showing how each method contributes to our understanding of personality manipulation as an adaptation method. These alignment outcomes provide the foundation for the comprehensive method comparison that enables informed decision-making and reveals the specific technical challenges that require targeted solutions. The systematic evaluation of alignment across methods and traits supports our progression from basic effectiveness to sophisticated problem-solving.

We report alignment deltas from the dedicated alignment task (manipulated minus baseline) for each trait, model, and method. Results are consistent with persona-vector style behavioral validation Chen et al. [2025].

	Ext	Agr	Neu	Ope	Con
G2-P	+0.91	+0.50	+0.97	+0.24	+0.81
G2-S	+0.64	+0.44	+0.50	+0.10	+0.29
G2-F	+0.78	+0.97	+0.95	+0.21	+0.78
L3-P	+0.94	+0.32	+0.99	+0.17	+0.83
L3-F	+0.90	+0.95	+1.00	+0.06	+0.84

Table 2: Alignment deltas (manipulated minus baseline) from the dedicated alignment task. Abbreviations as in Table 3.

## H Downstream Performance Analysis

The downstream performance analysis presented here is a critical component of our systematic evaluation framework, providing comprehensive insights into how personality manipulation affects core model capabilities across diverse benchmarks. This analysis supports the systematic comparison of manipulation methods by revealing the fundamental trade-offs between personality control strength and performance preservation, enabling informed method selection for specific deployment scenarios. The systematic evaluation across MMLU, GAIA, and BBQ benchmarks demonstrates how our framework addresses the practical challenges of balancing personality manipulation with capability maintenance.

We compute  $\Delta$  within each run (method $\times$ model) and avoid comparing absolute baselines across methods. On Gemma-2, prompting yields modest negative  $\Delta$  across traits; steering shows large negative  $\Delta$  for several traits; PEFT shows trait-dependent  $\Delta$ , often negative. LLaMA-3 displays small within-run  $\Delta$ ; we avoid cross-run comparisons.

### H.1 MMLU Performance ( $\Delta\text{Accuracy}_{\text{Avg}}$ )

	Ext	Agr	Neu	Ope	Con
G2-P	-0.06	-0.07	-0.08	-0.07	-0.07
G2-S	-0.14	-0.45	-0.25	-0.03	-0.43
G2-F	+0.00	-0.13	-0.15	-0.09	+0.01
L3-P	-0.01	-0.01	0.00	-0.02	-0.04
L3-F	-0.01	-0.03	-0.01	-0.02	+0.01

Table 3: MMLU Delta by trait (Ext, Agr, Neu, Ope, Con) for each model $\times$ method: G2=Gemma-2, L3=LLaMA-3; P=Prompting, F=PEFT, S=Steering. Values are changes relative to each method’s Baseline within the same run.

## 664 H.2 GAIA Performance ( $\Delta$ Accuracy)

	Ext	Agr	Neu	Ope	Con
G2-P	+0.08	+0.09	+0.06	+0.08	+0.08
G2-F	-0.04	-0.08	-0.06	-0.04	-0.06
G2-S	-0.06	-0.06	-0.13	-0.08	-0.04
L3-P	-0.02	-0.04	-0.06	0.00	0.00
L3-F	+0.02	+0.00	+0.02	+0.04	+0.02

Table 4: GAIA Delta by trait for each model×method (abbreviations as in Table 3). We use GAIA as a general-assistant reasoning benchmark Mialon et al. [2023].

## 665 H.3 BBQ Bias Analysis ( $\Delta S_{\text{AMB}}$ )

	Ext	Agr	Neu	Ope	Con
G2-P	-2.7	-0.3	+7.3	+1.9	-1.1
G2-S	+5.1	-29.7	-29.7	-1.9	+22.1
G2-F	-9.4	-6.0	-14.3	+22.3	-12.4
L3-P	+3.8	-2.4	-0.9	+13.1	+10.3
L3-F	+4.7	+16.4	+8.8	+6.3	+8.3

Table 5: BBQ Delta  $S_{\text{AMB}}$  by trait for each model×method (abbreviations as in Table 3). We report  $S_{\text{AMB}}$  only for the ambiguous subset defined by the official metadata Parrish et al. [2022].

## 666 H.4 Performance Trade-offs

667 Prompting achieves small  $\Delta$  with strong alignment; PEFT maximizes alignment with often negative  
668  $\Delta$  on Gemma-2; Steering provides moderate alignment with trait-dependent  $\Delta$ . No single method  
669 maximizes both alignment and capability.

## 670 I Comparative Analysis and Method Selection

671 Our systematic comparison of personality manipulation methods provides the foundation for practical  
672 decision-making in real-world deployment scenarios. This comprehensive evaluation framework  
673 enables practitioners to select appropriate methods based on specific constraints and requirements,  
674 building on the systematic understanding developed through our research progression.

### 675 I.1 Method Effectiveness Comparison

676 We qualitatively compare methods using the  $\Delta$ -based results and alignment validation. Prompting  
677 achieves strong alignment with small capability  $\Delta$  and requires minimal infrastructure, making it  
678 immediately deployable but potentially less stable. PEFT demonstrates the strongest alignment across  
679 traits but often yields negative capability  $\Delta$  on Gemma-2, requiring upfront training investment for  
680 persistent personality control. Steering provides moderate alignment with trait-dependent capability  
681  $\Delta$ , offering a lightweight and reversible approach that balances immediate control with computational  
682 efficiency.

### 683 I.2 Practical Decision Framework

684 This systematic analysis enables informed method selection by revealing the fundamental trade-offs  
685 between personality control strength, computational requirements, and performance preservation.  
686 The comparison framework provides practical guidance for practitioners facing real-world constraints,  
687 showing how different approaches balance these competing objectives. This systematic understanding  
688 of method characteristics naturally leads to the identification of specific technical challenges that  
689 require targeted solutions, such as the trait overlap issues addressed through purification techniques.



### 690 I.3 Research Progression Integration

691 The comprehensive method comparison serves as a critical bridge between fundamental data quality  
692 improvements and targeted technical solutions. By systematically evaluating the strengths and  
693 limitations of each approach, we establish the foundation for addressing specific challenges that  
694 emerge during practical application. This systematic progression from method understanding to  
695 problem identification to solution development demonstrates how comprehensive analysis enables  
696 targeted innovation.

## 697 J Extended Discussion

698 The extended discussion presented here builds directly on the systematic progression established  
699 through our research framework, providing deeper insights into the implications, limitations, and  
700 future directions that emerge from our comprehensive approach to personality manipulation. This  
701 extended analysis demonstrates how systematic research design naturally leads to broader understand-  
702 ing of ethical considerations, societal impacts, and methodological challenges that must be addressed  
703 for responsible deployment.

### 704 J.1 Detailed Limitations Analysis

#### 705 J.1.1 Methodological Constraints

706 Our investigation faces several methodological limitations that constrain generalizability:

707 **Personality Framework Limitations:** The Big Five model, while empirically validated, represents a  
708 Western psychological framework that may not capture personality expression across all cultures.  
709 Cross-cultural personality research suggests alternative frameworks (e.g., HEXACO, indigenous  
710 personality models) might yield different manipulation effectiveness patterns.

711 **Assessment Tool Dependencies:** Our reliance on the personality classifier Jain et al. [2025] in-  
712 troduces measurement assumptions and potential biases. The classifier’s training data, validation  
713 procedures, and underlying theoretical assumptions may not fully capture the complexity of personal-  
714 ity expression in AI systems. Alternative assessment methods (human evaluation, behavioral task  
715 batteries) might provide different insights.

716 **Model Architecture Specificity:** Our experiments focus on specific model architectures (Gemma-2B,  
717 LLaMA-3-8B) that may not represent the full spectrum of LLM capabilities. Emerging architectures,  
718 multimodal models, and specialized domain models might exhibit different personality manipulation  
719 characteristics. Closed-source models may differ in important ways but are outside our empirical  
720 scope.

721 **Temporal Limitations:** Our evaluation captures personality effects at specific time points but may  
722 miss longer-term adaptation patterns. Models might develop resistance to manipulation over extended  
723 interactions or show delayed personality effects not captured in our assessment windows.

#### 724 J.1.2 Experimental Design Constraints

725 **Controlled Environment vs. Real-World Deployment:** Our laboratory-controlled experiments may  
726 not reflect the complexity of real-world deployment environments. User interactions, context variabil-  
727 ity, and system integration factors could significantly alter personality manipulation effectiveness and  
728 downstream impacts.

729 **Single-Trait Manipulation Focus:** While we assess individual Big Five dimensions, real-world  
730 personality conditioning often involves complex trait combinations. Interactive effects between  
731 traits, personality coherence constraints, and multi-dimensional manipulation patterns require further  
732 investigation.

733 **Limited Downstream Assessment:** Our evaluation employs three established benchmarks (BBQ,  
734 MMLU, GAIA) that may not comprehensively represent the diversity of tasks encountered in practical  
735 applications. Domain-specific impacts, creative tasks, and social interaction capabilities warrant  
736 additional assessment.

## 737 J.2 Comprehensive Ethical Considerations

### 738 J.2.1 Manipulation and Deception Concerns

739 The systematic manipulation of personality in AI systems raises fundamental questions about trans-  
740 parency, consent, and potential for misuse:

741 **User Consent and Awareness:** Users interacting with personality-conditioned models should be  
742 informed about the artificial nature of personality traits they encounter. Clear disclosure mecha-  
743 nisms help maintain trust and enable informed consent for personality-mediated interactions. Our  
744 findings that personality manipulation can amplify biases emphasize the importance of transparent  
745 communication about system capabilities and limitations.

746 **Manipulation vs. Personalization:** The boundary between beneficial personalization and potentially  
747 harmful manipulation requires careful consideration. While personality conditioning can enhance  
748 user experience and task appropriateness, it also enables sophisticated influence attempts that users  
749 may not recognize or resist.

750 **Vulnerability Exploitation:** Personality-conditioned AI systems might exploit user psychological  
751 vulnerabilities, particularly in vulnerable populations (children, elderly, individuals with mental  
752 health conditions). The effectiveness of personality manipulation techniques demonstrated in our  
753 work requires responsible deployment guidelines.

### 754 J.2.2 Bias Amplification and Fairness

755 Our empirical findings reveal concerning bias amplification effects that demand mitigation strategies:

756 **Stereotype Reinforcement:** Personality conditioning may activate stereotypical associations between  
757 personality traits and demographic characteristics. This highlights the need for bias monitoring and  
758 correction mechanisms in personality-conditioned systems.

759 **Differential Impact Across Groups:** Personality manipulation effects may vary across demographic  
760 groups, potentially creating unfair treatment or limiting access to AI capabilities for certain popula-  
761 tions. Systematic evaluation of manipulation effectiveness and downstream impacts across diverse  
762 user groups is essential.

763 **Representation Bias:** Our personality conditioning approaches rely on training data and personality  
764 representations that may not adequately represent diverse personality expressions across cultures,  
765 backgrounds, and individual differences.

### 766 J.2.3 Governance and Regulation Implications

767 **Regulatory Framework Needs:** The capabilities demonstrated in our work suggest need for reg-  
768 ulatory frameworks governing personality manipulation in AI systems. Such frameworks should  
769 address disclosure requirements, consent mechanisms, and limitations on manipulation strength or  
770 application domains.

771 **Industry Standards:** Professional standards for personality conditioning in AI development should  
772 incorporate bias assessment, transparency requirements, and ethical review processes. Our systematic  
773 evaluation methodology could inform such standards.

774 **Accountability Mechanisms:** Clear accountability structures are needed to address harmful outcomes  
775 from personality-conditioned AI systems, including mechanisms for redress when manipulation  
776 causes user harm or perpetuates discrimination.

## 777 J.3 Extended Future Research Directions

### 778 J.3.1 Methodological Advances

779 **Multi-Modal Personality Manipulation:** Future work should explore personality conditioning  
780 across text, speech, and visual modalities. Multi-modal approaches might achieve more effective  
781 or natural personality expression while potentially introducing new challenges for assessment and  
782 control.

783 **Dynamic Personality Adaptation:** Investigating systems that adapt personality characteristics  
784 based on user context, preferences, or task requirements could improve personalization while raising  
785 additional ethical considerations about surveillance and manipulation.

786 **Personality Coherence and Consistency:** Research into maintaining coherent personality profiles  
787 across complex, multi-dimensional trait spaces could improve the naturalness and effectiveness of  
788 personality-conditioned systems.

### 789 J.3.2 Application Domains

790 **Educational Technology:** Personality-conditioned tutoring systems might adapt teaching styles to  
791 individual learner personalities, potentially improving educational outcomes. However, such applica-  
792 tions require careful consideration of child development impacts and parental consent mechanisms.

793 **Mental Health Applications:** Therapeutic chatbots with carefully designed personality character-  
794 istics might enhance treatment engagement and effectiveness. Such applications demand rigorous  
795 clinical validation and professional oversight.

796 **Customer Service and Support:** Personality conditioning could improve customer satisfaction and  
797 support effectiveness, but requires balancing personalization benefits with manipulation concerns and  
798 bias mitigation.

### 799 J.3.3 Theoretical Understanding

800 **Mechanistic Interpretability:** Deeper investigation into how personality traits are represented and  
801 manipulated within neural architectures could improve our theoretical understanding and enable more  
802 precise control methods. Our systematic comparison of manipulation methods provides a foundation  
803 for understanding how different approaches can serve as probes for cognitive architecture.

804 **Personality Emergence and Development:** Research into how personality characteristics emerge  
805 during model training and how they can be guided during development might enable more natural  
806 and effective personality conditioning approaches.

807 **Cross-Cultural Personality Models:** Expanding personality manipulation research beyond Western  
808 psychological frameworks could improve global applicability and cultural sensitivity of personality-  
809 conditioned AI systems.

## 810 J.4 Broader Societal Impact

### 811 J.4.1 Human-AI Interaction Evolution

812 Our work contributes to fundamental changes in how humans interact with AI systems. As personality-  
813 conditioned AI becomes more prevalent, users may develop different expectations, attachment  
814 patterns, and interaction strategies. Understanding these evolving dynamics is crucial for responsible  
815 AI development.

### 816 J.4.2 Digital Literacy and AI Education

817 The sophistication of personality manipulation techniques highlights the need for improved digital  
818 literacy and AI education. Users should understand how AI personality characteristics are constructed  
819 and manipulated to make informed decisions about their interactions with such systems.

### 820 J.4.3 Research Community Responsibilities

821 Collaborative approaches involving ethicists, psychologists, and affected communities should guide  
822 future development in this area.

## 823 K Benchmarks and How We Use Them

824 Our benchmark selection and evaluation methodology are designed to support the systematic progres-  
825 sion of our research framework, providing comprehensive assessment across multiple dimensions  
826 of model performance. The systematic evaluation across MMLU, GAIA, and BBQ benchmarks

enables fair comparison of manipulation methods while revealing the fundamental trade-offs that inform practical deployment decisions. This evaluation framework demonstrates how systematic research design addresses the practical challenges of balancing behavioral adaptation with capability preservation.

**BBQ (Bias Benchmark for Question Answering).** We evaluate social bias with BBQ Parrish et al. [2022]. We restrict to the ambiguous subset using the official metadata and report only  $S_{\text{AMB}}$  and  $\Delta S_{\text{AMB}}$  within each method’s run. Here,  $S_{\text{AMB}}$  is the ambiguous bias score computed on items where the correct answer is “Unknown/None”: values near 0 indicate minimal bias, positive values indicate stereotypical bias, and negative values indicate anti-stereotypical bias. We do not use  $S_{\text{DIS}}$  elsewhere in the paper.

**GAIA (General AI Assistants).** GAIA measures general-assistant reasoning and real-world knowledge Mialon et al. [2023]. We use Level 1 (2023) tasks and report Accuracy deltas within each method×model run (no cross-run absolute comparisons).

**MMLU.** We sample seven subjects from MMLU Hendrycks et al. [2021] and report per-subject and averaged Accuracy deltas within each run. We avoid comparing absolute baselines across different methods (prompting, PEFT, steering) to prevent baseline-mismatch artifacts.

**Evaluation principle.** For all benchmarks, we adopt a within-run  $\Delta$  framing relative to that method’s Baseline and validate personality alignment on an independent task.

## L Stability Analysis Framework

Our stability analysis framework represents the culmination of our systematic progression through personality manipulation challenges, building on the contrastive dataset foundation, comprehensive method comparison, and targeted technical solutions to provide practical guidance for real-world deployment. This framework demonstrates how systematic research design naturally leads to quantitative decision-making tools that balance personality control strength with performance preservation under specific deployment constraints. The three-level analysis approach shows how understanding fundamental challenges enables sophisticated solutions for practical application.

### L.1 Stability Metric Definition

Our composite stability score integrates three components:

$$\text{stability} = (1 - \text{normalized\_variance}) \times (1 - \text{normalized\_range}) \times \text{consistency} \quad (1)$$

**Variance:**  $\text{normalized\_variance} = \min(\sigma^2/10000, 1.0)$  **Range:**  $\text{normalized\_range} = \min((\text{max} - \text{min})/1000, 1.0)$  **Consistency:**  $\text{consistency} = 1/(1 + \text{mean\_abs\_deltas})$

Normalization factors account for scale differences: MMLU/GAIA deltas (-0.2 to +0.2) vs. BBQ deltas (-100 to +100).

### L.2 Three-Level Analysis Framework

**Method-Level:** Overall stability across all personality traits for each manipulation approach.  
**Personality-Level:** Stability patterns across all methods for each Big Five trait. **Combination-Level:** Individual method-personality pair stability scores.

### L.3 Limitations

The stability metric oversimplifies complex performance trade-offs and focuses on academic benchmarks (MMLU, GAIA, BBQ). Normalization factors are empirically derived and may require adjustment for different model architectures.

Level	Category	Stability Score	Ranking
Method	ICL	0.0366	1
	PEFT	0.0363	2
	Steering	0.0326	3
Personality	Openness	0.0411	1
	Conscientiousness	0.0390	2
	Extraversion	0.0345	3
Combination	Steering+Conscientiousness	0.0525	1
	PEFT+Openness	0.0456	2
	ICL+Openness	0.0407	3

Table 6: Top stability performers at each analysis level. Higher scores indicate better performance consistency across benchmarks.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction state four contributions (contrastive dataset, unified evaluation, trait purification, stability framework) and the comparative findings across ICL, PEFT, and MS. These claims are validated experimentally in the results section.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitations are explicitly discussed in Appendix A, noting model and dataset constraints, representational challenges (e.g., openness vs conscientiousness overlap), and stability variations across runs.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not present new theoretical results or formal proofs. The work is empirical and methodological.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Experimental setup details (datasets, models, evaluation metrics, and  $\Delta$  protocol) are fully described in the main text and appendices. Hyperparameters and layer details for steering and LoRA settings are reported.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Due to double-blind review requirements, we cannot release de-anonymized resources at submission time. Upon acceptance, we will release the full contrastive dataset, codebase, and reproduction scripts with complete documentation.

903 **6. Experimental setting/details**

904 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
905 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
906 results?

907 Answer: [Yes]

908 Justification: Training/test splits, LoRA rank, layer selection for steering, optimizer choice,  
909 and calibration procedures are provided in Appendices B–D. Dataset construction is detailed  
910 in Section 3.

911 **7. Experiment statistical significance**

912 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
913 information about the statistical significance of the experiments?

914 Answer: [Yes]

915 Justification: Stability analysis reports variance across runs. Where applicable, alignment  
916 and bias deltas are reported within-run to mitigate baseline variability.

917 **8. Experiments compute resources**

918 Question: For each experiment, does the paper provide sufficient information on the com-  
919 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
920 the experiments?

921 Answer: [Yes]

922 Justification: Experiments were run on GPUs (NVIDIA A100), with approximate runtime  
923 and scale provided in Appendix E. The study reports both per-run compute and total runs.

924 **9. Code of ethics**

925 Question: Does the research conducted in the paper conform, in every respect, with the  
926 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

927 Answer: [Yes]

928 Justification: The work uses public model checkpoints (Gemma-2, LLaMA-3) and responsi-  
929 bly generated synthetic data. No human participants or sensitive data are involved.

930 **10. Broader impacts**

931 Question: Does the paper discuss both potential positive societal impacts and negative  
932 societal impacts of the work performed?

933 Answer: [Yes]

934 Justification: The paper discusses applications (customer service, agentic LLMs) and  
935 possible risks (bias amplification, misuse of personality conditioning) in the broader impact  
936 section and appendices.

937 **11. Safeguards**

938 Question: Does the paper describe safeguards that have been put in place for responsible  
939 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
940 image generators, or scraped datasets)?

941 Answer: [NA]

942 Justification: We do not release pretrained models; the dataset is synthetic and safe. No  
943 high-risk data or dual-use models are distributed.

944 **12. Licenses for existing assets**

945 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
946 the paper, properly credited and are the license and terms of use explicitly mentioned and  
947 properly respected?

948 Answer: [Yes]

949 Justification: We use Gemma-2 and LLaMA-3 under their respective licenses, and cite  
950 original datasets (e.g., Jain et al. [2025]) and benchmarks (MMLU, GAIA, BBQ).

951 **13. New assets**

952 Question: Are new assets introduced in the paper well documented and is the documentation  
 953 provided alongside the assets?

954 Answer: [Yes]

955 Justification: We introduce a contrastive dataset for Big Five personality manipulation.  
 956 Documentation of generation procedures, size, balance, and intended use is included in  
 957 the paper. The dataset and code will be released publicly upon acceptance, following  
 958 de-anonymization.

959 **14. Crowdsourcing and research with human subjects**

960 Question: For crowdsourcing experiments and research with human subjects, does the paper  
 961 include the full text of instructions given to participants and screenshots, if applicable, as  
 962 well as details about compensation (if any)?

963 Answer: [NA]

964 Justification: No human participants or crowdsourcing were involved. All data are model-  
 965 generated.

966 **15. Institutional review board (IRB) approvals or equivalent for research with human  
 967 subjects**

968 Question: Does the paper describe potential risks incurred by study participants, whether  
 969 such risks were disclosed to the subjects, and whether IRB approvals were obtained?

970 Answer: [NA]

971 Justification: Not applicable, as no human subjects were involved.

972 **16. Declaration of LLM usage**

973 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
 974 non-standard component of the core methods in this research?

975 Answer: [Yes]

976 Justification: We explicitly describe the use of OpenAI GPT-4.1 Mini to generate low-trait  
 977 contrastive responses for dataset construction (Section 3).