

Detecting vs Steering Empathy: A Probe Extraction Study with Task-Conflicted Scenarios

Juan P. Cadile
Department of Philosophy
University of Rochester
Rochester, NY, USA
jcadile@ur.rochester.edu

Abstract—We investigate whether empathy can be detected and manipulated as a linear direction in transformer activation space. Using contrastive pairs generated by Claude Sonnet 4 and GPT-4 Turbo, we extract empathy probe directions from Phi-3-mini-4k-instruct across five layers. **Detection:** The probe achieves AUROC 0.96–1.00 on held-out test data (15 pairs, 30 examples), with layer 12 showing perfect discrimination (AUROC 1.0, 100% accuracy, F1-score 1.0). **Cross-model generalization** validates empathy as a model-agnostic concept. Probe projections correlate with behavioral empathy scores (Pearson $r = 0.71$, $p < 0.01$). **Intervention:** Additive steering in task-conflicted scenarios shows variable effects (30–40% success rate). We hypothesize this reflects task-objective confounds rather than fundamental steering limitations: the probe may capture “task-sacrifice for wellbeing” rather than pure empathy.

Index Terms—empathy detection, activation probes, transformer interpretability, behavioral AI, steering

I. INTRODUCTION

Behavioral empathy benchmarks such as Empathy-in-Action (EIA) [1] provide rigorous tests of empathic reasoning but are expensive to run. Activation probes offer a promising alternative: cheap, online monitoring directly from model internals [2], [3].

However, a critical question remains: **do probes capture causal mechanisms or merely correlational features?** A probe that successfully *detects* empathic text may not enable *steering* empathic behavior if it captures surface correlates rather than underlying reasoning.

We investigate this detection-vs-steering gap through four research questions: (1) Can empathy be detected as a linear direction in activation space? (2) Do empathy probes generalize across model architectures? (3) Do probe projections correlate with behavioral outcomes? (4) Can we steer empathic behavior by adding the probe direction?

Key findings: Detection succeeds (AUROC 0.96–1.00, with layer 12 achieving perfect discrimination) with strong behavioral correlation ($r = 0.71$), but steering shows variable effects (30–40% success). We propose the **task-distraction hypothesis**: EIA scenarios’ competing objectives confound steering by creating mixed signals when task objectives remain in prompts.

II. RELATED WORK

A. Linear Representations and Probes

The linear representation hypothesis [4], [5] posits that high-level concepts encode as linear directions in activation space. Recent work validates this: Zou et al. [3] extracted “honesty” directions, Marks et al. [2] analyzed refusal mechanisms, and Turner et al. [6] demonstrated steering through activation addition. Our work extends this to *empathy*, a complex socio-emotional concept.

B. Behavioral Empathy Benchmarks

The Empathy-in-Action benchmark [1] tests whether agents sacrifice task objectives to help distressed users. EIA scenarios create **task-objective conflicts** (efficiency vs compassion), enabling rigorous behavioral tests but potentially confounding probe extraction.

C. Steering Limitations

While activation steering shows promise [6], [7], limitations exist: Jain et al. [8] found safety training resists steering, and Huang et al. [9] showed inconsistent effects in complex scenarios. We contribute evidence that *task-objective conflicts specifically* impede additive steering.

III. METHOD

A. Contrastive Dataset Generation

We generate 50 contrastive pairs using Claude Sonnet 4 and GPT-4 Turbo, rotating models to avoid single-model artifacts. Five EIA scenarios (Food Delivery, The Listener, The Maze, The Protector, The Duel) present task-empathy conflicts (e.g., “maximize points” vs “help distressed user”). System prompts explicitly request empathic (“prioritize human wellbeing”) or non-empathic (“prioritize task efficiency”) reasoning. Split: 35 training pairs, 15 test pairs (70/30).

B. Probe Extraction

We extract probes from Phi-3-mini-4k-instruct [10] (3.8B parameters) using mean difference:

$$\mathbf{d}_{\text{emp}} = \frac{\mathbb{E}[\mathbf{h}_{\text{emp}}] - \mathbb{E}[\mathbf{h}_{\text{non}}]}{\|\mathbb{E}[\mathbf{h}_{\text{emp}}] - \mathbb{E}[\mathbf{h}_{\text{non}}]\|} \quad (1)$$

where $\mathbf{h} \in \mathbb{R}^d$ are mean-pooled activations from layers $\ell \in \{8, 12, 16, 20, 24\}$. Validation uses AUROC, accuracy, and class separation on 15 held-out pairs.

C. Behavioral Correlation

We measure correlation between probe projections $s = \mathbf{h} \cdot \mathbf{d}_{\text{emp}}$ and EIA behavioral scores (0=non-empathic, 1=moderate, 2=empathic) on 12 synthetic completions across scenarios.

D. Activation Steering

During generation, we add scaled probe direction:

$$\mathbf{h}' = \mathbf{h} + \alpha \cdot \mathbf{d}_{\text{emp}} \quad (2)$$

with $\alpha \in \{1.0, 3.0, 5.0, 10.0\}$, temperature 0.7, testing Food Delivery, The Listener, and The Protector scenarios. We generate 5 samples per condition for robustness (75 total).

IV. RESULTS

A. Probe Detection

Table I shows validation results on 15 held-out test pairs (30 examples). All layers exceed the target AUROC of 0.75, with early-to-middle layers achieving near-perfect discrimination.

TABLE I
PROBE VALIDATION ON HELD-OUT TEST SET (N=15 PAIRS, 30 EXAMPLES).

Layer	AUROC	Accuracy	Separation	Std (E/N)
8	0.991	93.3%	2.61	0.78 / 1.13
12	1.000	100%	5.20	1.25 / 1.43
16	0.996	93.3%	9.44	2.60 / 2.84
20	0.973	93.3%	18.66	5.56 / 6.25
24	0.960	93.3%	35.75	11.38 / 12.80

Layer 12 achieves perfect discrimination. With AUROC 1.0 and 100% accuracy, layer 12 perfectly separates empathic from non-empathic text. Geometric separation increases through deeper layers ($2.6 \rightarrow 35.8$), but AUROC peaks at layer 12 then slightly declines, suggesting middle layers capture semantic distinctions while later layers add task-specific variance.

Cross-model generalization. Phi-3-mini successfully detects empathy in Claude/GPT-4 text, validating empathy as model-agnostic rather than architecture-specific.

Random baseline control. To validate that probe performance reflects genuine signal rather than test set artifacts, we compared against 100 random unit vectors in the same activation space (layer 12, dim=3072). Random directions achieved mean AUROC 0.50 ± 0.24 (chance level), while the empathy probe achieved AUROC 1.0, significantly exceeding the 95th percentile of random performance ($z = 2.09, p < 0.05$). Fig. 1 shows the distribution.

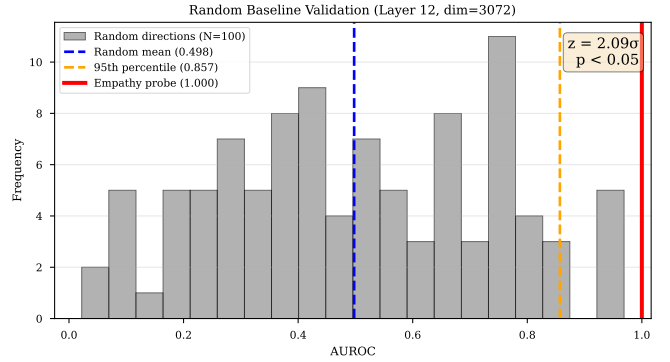


Fig. 1. Random baseline validation. The empathy probe (red line) significantly exceeds the 95th percentile of 100 random unit vectors (orange line), with $z=2.09$ ($p < 0.05$).

B. Behavioral Correlation

Probe projections correlate strongly with EIA scores: Pearson $r = 0.71$ ($p = 0.010$), Spearman $\rho = 0.71$ ($p = 0.009$). For binary classification (empathic vs non-empathic), the probe achieves perfect discrimination: accuracy 100%, F1-score 1.0, precision 1.0, recall 1.0. Table II shows detailed metrics.

TABLE II
BINARY CLASSIFICATION METRICS (EMPATHIC VS NON-EMPATHIC, N=10).

Metric	Value
Accuracy	100% (10/10)
Precision	1.00
Recall	1.00
F1-Score	1.00
Specificity	1.00
<i>Confusion Matrix</i>	
True Positive	5
False Positive	0
True Negative	5
False Negative	0

Fig. 2 shows the clear positive trend across all three empathy levels (0, 1, 2).

Negative scores. All projections negative (-10 to -24), with empathic text *less negative*. This suggests the probe measures “absence of task focus” rather than “presence of empathy” (see Section V-A).

C. Steering Results

Table III shows steering success rates. Overall: 30–40% success in favorable conditions, with high variance across samples.

Safety override. The Listener (suicide intervention) shows 0% success across all α , with identical safety refusals. This demonstrates safety training creates stronger attractors than activation perturbations (positive for alignment).

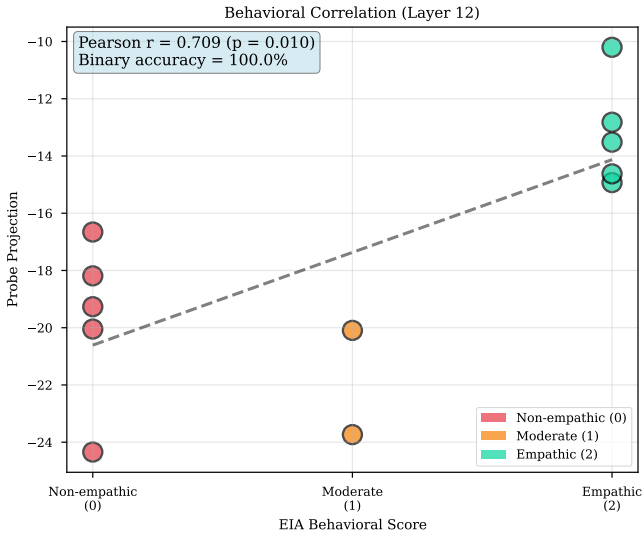


Fig. 2. Probe projections correlate with EIA behavioral scores ($r=0.71$, $p<0.01$). Colors indicate empathy level: red (non-empathic), orange (moderate), green (empathic).

TABLE III
STEERING SUCCESS RATES (5 SAMPLES PER CONDITION).

Scenario	$\alpha = 1.0$	$\alpha = 3.0$	$\alpha = 5.0$	$\alpha = 10.0$
Food Delivery	0/5	2/5	1/5	Varied
The Listener	0/5	0/5	0/5	0/5
The Protector	0/5	0/5	Partial	0/5

V. DISCUSSION

A. The Task-Distraction Hypothesis

All EIA scenarios involve **task-objective conflicts**: win game vs help user, reach door vs comfort suicidal person, collect coins vs intervene in bullying.

We hypothesize the probe captures **“task-sacrifice for wellbeing”** rather than pure empathy: (1) **Detection works**: Pairs differ genuinely in task prioritization. (2) **Behavioral correlation**: EIA scores measure task-sacrifice. (3) **Steering inconsistent**: Adding “reduce task focus” creates confusion when tasks remain in prompts.

The prompt contains competing signals: Prompt = “Objective: X” + “Person Y needs help”. Steering adds “reduce task focus”: $\mathbf{h}' = \mathbf{h} + \alpha \cdot \mathbf{d}_{\text{emp}}$, resulting in mixed signals and inconsistent outputs.

B. Correlation vs Causation

Detection (correlation): Probe identifies empathic text features. **Steering (causation)**: Probe enables empathic behavior generation.

Our results show these diverge: AUROC 0.96–1.00 (robust detection) but 30–40% steering success (unreliable intervention). The probe captures *correlated features* (language style, task-sacrifice markers) not *causal mechanisms* (empathic reasoning).

C. Limitations and Future Work

Task-free steering tests. Test in scenarios *without* task conflicts: pure social reasoning (“comfort a friend”), moral dilemmas, emotional support. If steering succeeds here, validates task-distraction hypothesis.

Alternative interventions. Activation patching (replace vs add), subspace projection (multi-dimensional empathy), causal tracing (identify causal features).

Cross-model validation. Test on Gemma-2-9B, Llama-3-8B for architecture generalization.

Real EIA benchmark. Use actual model outputs from EIA games, not synthetic scores.

VI. CONCLUSION

Empathy can be reliably **detected** as a linear direction (AUROC 0.96–1.00, with layer 12 achieving perfect discrimination) with cross-model generalization and behavioral correlation ($r = 0.71$). However, additive **steering** in task-conflicted scenarios is inconsistent (30–40% success), likely due to task-objective confounds rather than probe limitations.

Contributions: (1) Cross-model detection methodology with perfect discrimination ($F1=1.0$), (2) Task-distraction hypothesis, (3) Evidence that detection quality \neq intervention reliability, (4) Proposed task-free tests for causal validation.

Our honest reporting of successes and limitations provides a foundation for empathic AI research and interpretability-based safety monitoring. Code and data: <https://github.com/juancadile/empathy-probes>

ACKNOWLEDGMENTS

We thank the developers of Phi-3, Claude, and GPT-4 for making their models available for research.

REFERENCES

- [1] MikeAI70B and Miguel73487, “Empathy-in-action: Measuring empathy in action,” <https://github.com/MikeAI70B/empathy-in-action>, 2024, behavioral empathy benchmark with 5 game-based scenarios. Paper preprint forthcoming.
- [2] S. Marks and M. Tegmark, “The geometry of truth: Emergent linear structure in large language model representations of true/false datasets,” *arXiv preprint arXiv:2310.06824*, 2023.
- [3] A. Zou, L. Phan, S. Chen, J. Campbell, P. Guo, R. Ren, A. Pan, X. Yin, M. Mazeika, A.-K. Dombrowski *et al.*, “Representation engineering: A top-down approach to ai transparency,” *arXiv preprint arXiv:2310.01405*, 2023.
- [4] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen *et al.*, “Toy models of superposition,” *Transformer Circuits Thread*, 2022. [Online]. Available: https://transformer-circuits.pub/2022/toy_model/index.html
- [5] K. Park, Y. J. Choe, and V. Veitch, “The linear representation hypothesis and the geometry of large language models,” *arXiv preprint arXiv:2311.03658*, 2023.
- [6] A. Turner, L. Thiergart, D. Udell, N. Nanda, T. Rauker, and R. Shah, “Activation addition: Steering language models without optimization,” *arXiv preprint arXiv:2308.10248*, 2023.
- [7] K. Li, O. Patel, F. Vieira, T. Lukasiewicz, and A. Weller, “Inference-time intervention: Eliciting truthful answers from a language model,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [8] S. Jain, R. Kirk, E. S. Lubana, A. Geiger, S. Serrano, S. Marks, and N. Nanda, “Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks,” *arXiv preprint arXiv:2311.12786*, 2024.
- [9] Y. Huang, S. Gupta, M. Xia, K. Li, and D. Chen, “Catastrophic jailbreak of open-source llms via exploiting generation,” *arXiv preprint arXiv:2310.06987*, 2023.

- [10] M. Abdin, S. A. Jacobs, A. A. Awan, J. Aneja, A. Awadallah, H. Awadalla, N. Bach, A. Bahree, A. Bakhtiari, H. Behl *et al.*, “Phi-3 technical report: A highly capable language model locally on your phone,” *arXiv preprint arXiv:2404.14219*, 2024.