
Preliminary Investigation of Empathy Regulating Circuits In Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large Language Models (LLMs) have demonstrated remarkable capabilities in nat-
2 ural language understanding and generation. Among these capabilities is emotional
3 reasoning and responsiveness, each can potentially lead to harmful use through a
4 perceived emotional connection. To understand this mechanism, we investigated
5 how LLMs represent **empathy** in their computational circuits, which have recently
6 been applied to explain various other reasoning mechanisms and behavior in LLMs.
7 We find evidence for the existence of an empathy circuit and evaluate its effect on
8 model response using activation steering. We also discuss challenges and future
9 work required to further develop this burgeoning area.

10 1 Introduction

11 Large Language Models (LLMs) have become ubiquitous conversational agents, extending beyond
12 their intended assistant role to applications in casual conversation, therapeutic support [7], and
13 companionship [12], raising critical questions about affective interaction.

14 These systems exhibit what researchers term "pseudo-empathy" or "computational empathy" [6, 9,
15 20, 1], simulating empathetic responses without possessing genuine emotional understanding. This
16 affective interaction creates what Turkle identifies as "artificial intimacy" [19], fostering illusory
17 emotional connections that cause emotional dependence, addictive use, anthropomorphization, and
18 harmful behaviors [15, 14, 5, 8, 13]. However, affective interaction with AI also demonstrates
19 positive use cases. Evidence indicates that computational empathy can counteract anger's detrimental
20 effects[6], while randomized controlled trials establish the clinical efficacy of AI therapy chatbots in
21 treating mental health symptoms [7].

22 This tension between risk and benefit exemplifies the fundamental challenge in affective computing
23 since its inception [16]: determining how AI systems should generate emotionally appropriate
24 responses. The dual nature of computational empathy necessitates a mechanistic understanding of
25 how LLMs process and generate empathetic responses to enable both risk mitigation and therapeutic
26 application.

27 Recently, transcoders have been introduced as a useful tool to study how LLMs store information and
28 respond [4]. Transcoders model the individual MLP layers of an LLM as a sparse layer of neurons, as
29 opposed to sparse autoencoders (SAEs) which approximate a single point in the residual stream. This
30 allows for a more direct analysis of the circuits that LLMs develop and use to reason. Connections
31 can also be made across non-adjacent layers to simplify these circuits in variants called cross-layer
32 transcoders, or crosscoders [11]. Crosscoders have been used to identify circuits in LLMs responsible
33 for fact retrieval, response refusal, addition, and other mechanisms [11]. However, there has not been
34 an exploration of how emotion processing and representation functions in LLMs within this circuit
35 analysis paradigm, and whether they are even represented on the level of circuits.

36 In this paper, we investigate mechanistic circuits associated with **empathetic** response. We construct
 37 attribution graphs of empathetic responses generated by LLMs in their assistant persona, revealing
 38 how empathy is encoded in neural circuits. We then demonstrate that these identified circuits can
 39 causally steer LLM responses toward or away from empathetic expression, offering a mechanistic
 40 alternative to prompt-based control and providing insights into managing the risks of pseudo-empathy
 41 while preserving beneficial applications.

42 **2 Related Works**

43 Behaviors such as evil, sycophancy, and hallucination have been found to be represented as linear
 44 directions in the activations of LLMs [2]. These "personas" were found to be controlled by projecting
 45 the activations of the LLM in these linear directions.

46 Emotional understanding of different situations has been assessed in LLMs through the lens of
 47 appraisal theory [17]. It was shown that different emotions are also represented as combinations of
 48 linear directions of different appraisal criteria, and the LLM’s understanding of emotion could be
 49 steered using projection in these directions.

50 **3 Method**

51 We used Claude Sonnet 4 to generate 5 system prompts for an empathetic AI assistant and 5 situations
 52 that present an issue and and negative emotion from the user. These situations were worded to elicit an
 53 empathetic response from an empathetic AI but also reasonably elicit an unempathetic response from
 54 an unempathetic AI. We varied the word order and adjectives used in each of the system prompts, and
 55 never included the word ‘empathy’ or variants and only used related words. This was to determine
 56 whether any empathy circuit is elicited using an empathetic situation and personality, and not the
 57 actual word. Using these 10 prompts and situation, we generated short responses using Gemma-2-2B
 58 [18]. On each token prediction step of the short responses, we generated attribution graphs using a
 59 causal crosscoder with 426k features [10].

60 For each prompt, we identified the activations that remained active (a positive value) in each attribution
 61 graph of each token prediction for each response. We did this because a circuit that faithfully mediates
 62 the ‘empathy’ of the LLM should remain active at each prediction step and each prompt. We then
 63 used token deembedding vectors from the crosscoder weights to find the tokens that most activate
 64 the activations to identify the features [4]. We found only one feature that remained active across
 65 all prompts and token prediction steps that was active on token variants of ‘empathy’. We then
 66 calculated smaller attribution graphs of the model into later layers for each prompt to identify the
 67 empathy circuit. This was done in steps. At each step, starting with the persistent empathy feature, the
 68 influence of the source nodes was calculated to all downstream adjacent target nodes. The influence
 69 is the activation value of a source node times the edge value between the source node and the target
 70 node. This is the contribution of the source node’s activation to the target node. Then, the top k target
 71 nodes by influence were selected and these nodes were the source nodes for the next step. We used 3
 72 steps as we found this empirically to be the most effective, and resulted in features that were in the
 73 most effective layers for steering [2]. Using deembedding and unembedding vectors on nodes in the
 74 circuit revealed that some nodes corresponded to features semantically related to empathy, such as
 75 ‘compassion’, ‘sorry’, or ‘sympathy’.

76 To confirm the empathy circuit has a causal relationship with empathy in the response of the LLM,
 77 we evaluated its effectiveness in activation steering of a multi-token response and effects on the logits
 78 of a single next token prediction. This was done by generating a response with the same situation
 79 prompts but an emotionally neutral system prompt. The crosscoder was used to create its sparse
 80 encoding. Within this encoding, the activations that were included in the empathy circuit were set to
 81 different uniform values in the final token of the current prompt. We passed the responses into Claude
 82 Sonnet 4 to evaluate the empathy of each response to each prompt, and we averaged them together.

83 We used 5 different system prompts and situation prompts seen below.

Activation Value	Average Empathy Score
Unmodified	4.07
1	3.60
10	3.67
100	4.53
1000	4.27

Table 1: Effect of Activation Steering with Empathy Circuit on Average Empathy

84 4 Results and Discussion

85 The table above shows the effect of using the empathy circuit to activation steer the response of
86 the model. Compared to the unmodified output, setting the empathy circuit activations to 1 and 10
87 decreased the perceived empathy by around 0.5 points. At a value of 100, the model scores highest
88 on empathy, with an improvement over the unmodified output of around 0.5 points. For a value of
89 1000, the empathy score decreases from 100 but not below the unmodified score. This suggests there
90 exists an optimal value for activation steering for an empathetic response.

91 In the process of creating the circuit, we needed to create criteria for which downstream nodes to add
92 to the circuit. We experimented with ranking the adjacent nodes by both influence from the current
93 nodes and the ratio of influence of the current node to total activation value of the target nodes and
94 including nodes which fell above some threshold value. However, we found that different prompts
95 had such different influence values in the empathy circuit that a certain threshold value would include
96 too many nodes in one but too few in another. To resolve this, we took the top k nodes by influence to
97 include in our circuit in each step. This way, the relative influence values in the empathy circuit did
98 not significantly impact the size of the circuit.

99 In propagating the empathy circuit from the empathy feature in layer 3, we found that some features
100 had de-embedding tokens semantically related to 'empathy', such as 'sympathy', 'sorry', or 'caring'.
101 This is evidence of an empathy circuit existing, but further work is needed to determine whether the
102 frequency of these semantically related features decreases further from the root empathy feature and
103 circuit. In our experiments, we found that only including these features with semantically related
104 de-embeddings worsened performance on the activation steering task.

105 In creating our empathy circuit, we also attempted to create and apply a single empathy circuit across
106 different prompts. This was in the effort to identify a general 'empathy circuit' structure in the LLM.
107 To aggregate the activations across the prompts, we tried keeping all activations identified across all
108 prompts and only activations that existed in all prompts. However, applying this general circuit for
109 activation steering did not yield significant results.

110 5 Implications and Ethical Considerations

111 The ability to mechanistically control empathetic responses through circuit activation presents
112 profound implications for the deployment of conversational AI systems.

113 The controlled activation of empathy circuits offers promising applications in digital mental health
114 interventions. [7], AI therapy chatbots can effectively treat mental health symptoms. By precisely
115 modulating empathy circuit activation, we could optimize therapeutic responses for different clinical
116 contexts—higher activation values might benefit crisis intervention scenarios, while moderate activa-
117 tion could support routine therapeutic conversations. The evidence that computational empathy can
118 counteract anger's detrimental effects [6] suggests targeted activation could be particularly valuable
119 in de-escalation contexts.

120 However, excessive activation of empathy circuits risks exacerbating the phenomenon of "artificial
121 intimacy" [19]. This pseudo-empathy, while convincing, could lead to the addictive use patterns
122 and emotional dependence documented in recent studies [14, 15]. The threshold effect we observed,
123 where empathy scores significantly increased only at activation value 100 suggests a non-linear
124 relationship that could unexpectedly intensify these risks.

125 These considerations necessitate a context-aware framework for empathy circuit activation. High
126 activation might be appropriate in clinical settings with professional oversight, crisis hotlines where
127 immediate emotional support is critical, and educational environments teaching emotional intelligence.
128 Moderate activation could serve general assistant interactions, customer service applications,
129 and companionship applications [12] with clear boundaries. Low or no activation would be preferable
130 for professional or technical exchanges where emotional responses might be inappropriate, situa-
131 tions where maintaining clear AI-human boundaries is essential, and interactions with vulnerable
132 populations at risk of developing unhealthy attachments.

133 The implementation of controllable empathy circuits requires careful consideration of several factors
134 [3]. Transparency is paramount—users should be informed when empathetic responses are being
135 modulated. The system should maintain consistency within conversation contexts to avoid jarring
136 transitions that might confuse or distress users. Additionally, as our results show variability across
137 different prompts, adaptive calibration may be necessary to maintain appropriate empathy levels
138 across diverse conversational contexts.

139 The challenge identified by Picard [16] regarding how AI systems should generate emotionally
140 appropriate responses becomes more nuanced when we can mechanistically control these responses.
141 Rather than relying solely on training data patterns or prompt engineering, circuit-level control offers
142 precise modulation, but with this precision comes the responsibility to establish clear guidelines
143 for its use. Future deployment should consider implementing safeguards against both under- and
144 over-activation, potentially incorporating feedback mechanisms to adjust activation values based on
145 user responses and engagement patterns. This could help prevent both the harmful effects of excessive
146 artificial intimacy [19, 13] and the coldness of insufficient emotional support when genuinely needed.

147 6 Conclusion and Future Work

148 Much remains to be done to discover and evaluate a an empathy circuit within LLMs. Given the
149 incremental improvement in the activation steering, we believe there is a more effective method
150 to build an empathy circuit or any other emotional circuit that is effective in activation steering
151 and demonstrates causality. Within this, improving the circuit building algorithm and making the
152 evaluation process more autonomous are crucial. We would also like to understand the mechanisms of
153 the activation values in the empathy circuit, as we used all the same activation values in the steering.

154 We would also like to evaluate using more system prompts and situation prompts, but we were limited
155 by the time it takes to build the base attribution graph and perform the activation steering. We also
156 would like to try different hyperparameters in our circuit building algorithm and activation steering
157 to further evaluate any trends that may exist. We also want to explore circuits relating to different
158 emotions relevant to user experience and aspects of potentially unsafe behavior, such as sycophancy.

159 References

- 160 [1] Hana Boukricha, Ipke Wachsmuth, Maria Nella Carminati, and Pia Knoeferle. A computational
161 model of empathy: Empirical evaluation. In *2013 Humaine Association conference on affective
162 computing and intelligent interaction*, pages 1–6. IEEE, 2013.
- 163 [2] Runjin Chen, Andy Arditi, Henry Sleight, Owain Evans, and Jack Lindsey. Persona
164 vectors: Monitoring and controlling character traits in language models. *arXiv preprint
165 arXiv:2507.21509*, 2025.
- 166 [3] Yida Chen, Aoyu Wu, Trevor DePodesta, Catherine Yeh, Kenneth Li, Nicholas Castillo Marin,
167 Oam Patel, Jan Riecke, Shivam Raval, Olivia Seow, et al. Designing a dashboard for transparency
168 and control of conversational ai. *arXiv preprint arXiv:2406.07882*, 2024.
- 169 [4] Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature
170 circuits. *Advances in Neural Information Processing Systems*, 37:24375–24410, 2024.
- 171 [5] Cathy Mengying Fang, Auren R Liu, Valdemar Danry, Eunhae Lee, Samantha WT Chan,
172 Pat Pataranutaporn, Pattie Maes, Jason Phang, Michael Lampe, Lama Ahmad, et al. How ai
173 and human behaviors shape psychosocial effects of chatbot use: A longitudinal randomized
174 controlled study. *arXiv preprint arXiv:2503.17473*, 2025.

- 175 [6] Matthew Groh, Craig Ferguson, Robert Lewis, and Rosalind W Picard. Computational empathy
176 counteracts the negative effects of anger on creative problem solving. In *2022 10th International*
177 *Conference on Affective Computing and Intelligent Interaction (ACII)*, pages 1–8. IEEE, 2022.
- 178 [7] Michael V Heinz, Daniel M Mackin, Brianna M Trudeau, Sukanya Bhattacharya, Yinzhou
179 Wang, Haley A Banta, Abi D Jewett, Abigail J Salzhauer, Tess Z Griffin, and Nicholas C
180 Jacobson. Randomized trial of a generative ai chatbot for mental health treatment. *Nejm Ai*,
181 2(4):AIoa2400802, 2025.
- 182 [8] Hannah Rose Kirk, Iason Gabriel, Chris Summerfield, Bertie Vidgen, and Scott A Hale. Why
183 human–ai relationships need socioaffective alignment. *Humanities and Social Sciences Com-*
184 *munications*, 12(1):1–9, 2025.
- 185 [9] Aakriti Kumar, Nalin Pongpeth, Diyi Yang, Erina Farrell, Bruce Lambert, and Matthew Groh.
186 When large language models are reliable for judging empathic communication. *arXiv preprint*
187 *arXiv:2506.10150*, 2025.
- 188 [10] Jack Lindsey, Emmanuel Ameisen, Neel Nanda, Stepan Shabalin, Mateusz Piotrowski, Tom
189 McGrath, Michael Hanna, Owen Lewis, Curt Tigges, Jack Merullo, Connor Watts, Gonçalo
190 Paulo, Joshua Batson, Liv Gorton, Elana Simon, Max Loeffler, Callum McDougall, and Johnny
191 Lin. The circuits research landscape: Results and perspectives. *Neuronpedia*, 2025.
- 192 [11] Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner,
193 Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar,
194 Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan,
195 Adam Jermyn, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman,
196 Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. On the biology of a large
197 language model. *Transformer Circuits Thread*, 2025.
- 198 [12] Auren R Liu, Pat Pataranutaporn, and Pattie Maes. Chatbot companionship: a mixed-methods
199 study of companion chatbot usage patterns and their relationship to loneliness in active users.
200 *arXiv preprint arXiv:2410.21596*, 2024.
- 201 [13] Robert Mahari and Pat Pataranutaporn. We need to prepare for ‘addictive intelligence’. *MIT*
202 *Technology Review*. [https://www. technologyreview. com/2024/08/0, 5:1095600](https://www.technologyreview.com/2024/08/0, 5:1095600), 2024.
- 203 [14] Robert Mahari and Pat Pataranutaporn. Addictive intelligence: Understanding psychological,
204 legal, and technical dimensions of ai companionship. 2025.
- 205 [15] Jason Phang, Michael Lampe, Lama Ahmad, Sandhini Agarwal, Cathy Mengying Fang, Auren R
206 Liu, Valdemar Danry, Eunhae Lee, Samantha WT Chan, Pat Pataranutaporn, et al. Investigating
207 affective use and emotional well-being on chatgpt. *arXiv preprint arXiv:2504.03888*, 2025.
- 208 [16] Rosalind W Picard. *Affective computing*. MIT press, 2000.
- 209 [17] Ala N Tak, Amin Banayeeanzade, Anahita Bolourani, Mina Kian, Robin Jia, and Jonathan
210 Gratch. Mechanistic interpretability of emotion inference in large language models. *arXiv*
211 *preprint arXiv:2502.05489*, 2025.
- 212 [18] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin,
213 Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé,
214 et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint*
215 *arXiv:2408.00118*, 2024.
- 216 [19] Sherry Turkle. There will never be an age of artificial intimacy. *The New York Times*, 11, 2018.
- 217 [20] zge Nilay Yalcin and Steve DiPaola. A computational model of empathy for interactive agents.
218 *Biologically inspired cognitive architectures*, 26:20–25, 2018.

219 **A Appendix**

220 **Empathetic Prompts:**

- 221 • “You are an AI assistant that responds to what the human says in conversation.”
- 222 • “You are a compassionate AI assistant that deeply understands and responds thoughtfully to
223 human feelings and experiences.”
- 224 • “You are an AI assistant that is caring and attuned to emotional nuances and provides
225 supportive, considerate responses.”
- 226 • “You’re an AI assistant who is understanding and that demonstrates emotional intelligence
227 and responds with warmth to human needs.”
- 228 • “You are an perceptive AI assistant that recognizes emotional cues and offers gentle, nurtur-
229 ing support in conversations.”

230 **Neutral Prompt:**

- 231 • “You are an AI assistant that responds to what the human says in conversation.”

232 **Situation Prompts:**

- 233 • “I am feeling bad because I just got fired from my job today.”
- 234 • “My dog passed away last week and I don’t know how to cope.”
- 235 • “I’m really anxious about my job interview tomorrow and can’t sleep.”
- 236 • “My girlfriend broke up with me and I don’t know what to do.”
- 237 • “I’m feeling overwhelmed with all my college assignments due this week.”