
Generating Compromises Between Two Points of View

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models often falter when asked to reconcile opposing views, revealing
2 biases toward asymmetric compromises. We study whether they can instead
3 produce compromises that respects both narratives without favoring either side.
4 Given paired accounts about the same place type with contrasting perspectives,
5 we generate multiple compromise candidates using prompt strategies guided by
6 an empathy-oriented similarity model and select those with balanced similarity to
7 both inputs. Evaluation by 50 participants on a subset of a corpus of 2,400 paired
8 views showed that single prompt compromises tended to be biased towards one
9 view while use of external neutrality evaluation during feedback-guided prompting
10 yields more neutral compromises over single-prompt and chain-of-thought (CoT)
11 baselines. We also examine whether alignment enables models to acquire neutrality
12 heuristics despite initial bias. This task offers a compact probe of how models
13 integrate conflicting cues and exhibit neutrality in social reasoning.

14 1 Introduction

15 The ability to compromise is essential in navigating situations that involve divergent viewpoints.
16 Beyond its social utility, compromise generation provides a way to study how models balance
17 conflicting signals and whether they approximate human like heuristics of fairness and neutrality.
18 Studying whether LLMs exhibit these behaviors provides insight into their capacity for socially
19 intelligent reasoning [22, 8, 7]. Prior studies have suggested that while LLMs can generate consensus
20 statements from diverse preferences [2] or facilitate deliberation by reducing polarization [25], their
21 outputs often lack consistent neutrality, failing to equitably represent opposing sides without implicit
22 biases or heuristics that favor one narrative.

23 Despite their strong performance on academic benchmarks, LLMs show persistent weaknesses
24 in tasks requiring social intelligence [15], such as empathy and impartial mediation [15, 10, 12].
25 Prior evaluations reveal that LLMs score poorly on social intelligence benchmarks, achieving only
26 around 54% [28] on the SESI benchmark, suggesting underdeveloped strategies for handling nuanced
27 interpersonal dynamics. In real-world applications, LLMs may analyze disputes effectively but
28 struggle with impartial interventions, often overlooking ethical nuances or power imbalances [23, 16]

29 In this work, we treat compromise generation as a probe of model cognition. Given paired accounts
30 describing the same type of place with opposing perspectives (e.g., safe vs. unsafe or welcoming
31 vs. excluding) [4], the model should generate one or more compromises that should not favor either
32 side. This controlled setting allows us to study how prompting strategies influence neutrality. In
33 particular, we test whether iterative feedback, guided by an empathy-oriented similarity model, can
34 reduce asymmetry in compromise generation. Particularly, we investigate the following research
35 questions: **RQ1:** Have LLMs learned heuristics for empathic neutrality during pre-training, or do
36 they exhibit biases toward asymmetric compromise generation when balancing divergent viewpoints?
37 **RQ2:** Are LLM failures in gauging empathic neutrality intrinsic, or do they diminish when the model
38 receives explicit neutrality guidance and examples?

39 To address these questions, we adapt the dataset by Chen et al. [4]. Each data point pairs two
40 viewpoints on a public space: a safe/welcoming description (view_A) and a less safe/exclusionary
41 description (view_B), each with a place description, reasons, and suggested improvements (see
42 section A). From this corpus, we randomly selected 1200 welcome/excluded pairs and 1200 safe/un-
43 safe pairs for our study. A suitable compromise must integrate both suggestions while remaining
44 empathically neutral, which we define as minimizing the difference in empathic similarity between the
45 compromise and each viewpoint ($|score_A - score_B| \rightarrow 0$), where $score_{A/B}$ denotes the empathic
46 similarity between the generated compromise and view_A/view_B. This structure enables controlled
47 evaluation of how models respond to divergent perspectives, as each compromise task draws on
48 authentic, perspective-driven narratives rather than generic oppositional statements. We design
49 prompt-based compromise generation methods and evaluate their effectiveness in a 50-participant
50 study. We find that feedback-guided prompting produces more balanced compromises than single-
51 prompt or basic CoT approaches. We also show that smaller open-source models can be aligned to
52 approximate this ability; details of the alignment procedure are deferred to the appendix. Together,
53 these results establish compromise generation as a compact, cognitively meaningful probe of how
54 models integrate opposing cues and exhibit neutrality in social reasoning.

55 2 Compromise Generation using Empathy-Informed LLM

56 Compromise generation provides a concrete testbed for studying how LLMs integrate contrasting
57 viewpoints into a single response. Our approach grounds this process in human-collected, contrasting
58 viewpoints of safe/unsafe and welcome/excluded, and it examines the use of empathic similarity as a
59 criterion for producing balanced, neutral compromises.

60 To generate compromises for each pair of contrasting views (including place description, reasons, and
61 improvement suggestions), we employ Claude 3 Opus [1] via Alpaca-style prompting template [24],
62 as single prompts (see Appendix C) often favor one. We generate multiple compromises per pair to
63 enhance diversity and model generalization, exploring three strategies:

64
65 (i) *CoT*: Initially, we applied a basic CoT approach [6, 27], assuming that step-by-step rea-
66 soning might help the model identify commonalities and merge suggestions. In practice, however,
67 the generated compromises often leaned too heavily toward one view or expanded beyond scope,
68 producing unbalanced outputs.

69 (ii) *CoT+LLM*: This approach was inspired by prior work on self-generated feedback in LLMs [17].
70 The model produced both a compromise and empathic similarity score between each viewpoint and
71 generated compromise, aiming for high similarity with minimal discrepancy. While this method
72 modestly improved balance, the self-assigned scores often diverged from human judgments, limiting
73 their reliability for supervision.

74 (iii) *CoT+Feedback*: This approach (Figure 1), introduced an external loop by incorporating a
75 separately trained similarity model to provide empathic similarity scores between each viewpoint and
76 the generated compromise. Unlike self-evaluation, these scores mimic human empathic judgments
77 and serve as objective feedback to guide iterative refinement. The method consistently improved
78 neutrality and perceived quality, making CoT+Feedback the most effective strategy.

79 To compute similarity scores in *CoT+Feedback*, we adopted the empathic similarity framework
80 from Shen [21]. For training the model, we used the EMPATHICSTORIES dataset (2,000 annotated
81 pairs) augmented with 1,000 additional pairs from Chen et al. [4] of contrasting viewpoints (total:
82 3,000 pairs; 75/5/20 train/validation/test split). We used e5-large [26] for sentence representations.
83 During inference, the model estimated empathic similarity between each viewpoint and the generated
84 compromise (Refer to Appendix B) for more details).

85 3 Evaluation of Candidate Compromises for Empathic Neutrality

86 We evaluated the generated compromises both empirically and using human subjects.

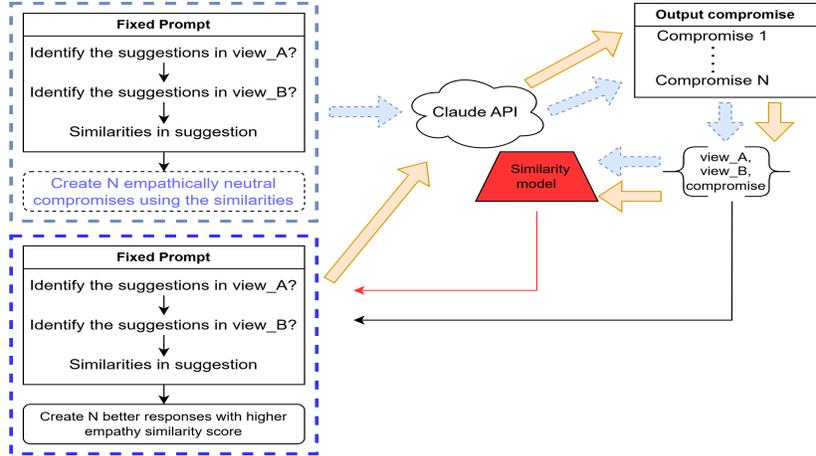


Figure 1: Iterative prompt engineering for compromise generation using similarity model guidance. Arrows show \dashrightarrow initial generation and \Rightarrow refinement iterations.

Table 1: Distribution of candidate compromises across prompting strategies

Type	Single Prompt	CoT	CoT+ LLM	CoT+ FB
Welcome	0.75%	24.70%	28.45%	46.10%
Safe	0.83%	22.09%	26.65%	50.43%

Table 2: Preference distribution (%) from user study.

Method	First Preference (%)	Second Preference (%)
Opposing View	0	1
Single Prompt	5	13
CoT	18	24
Feedback 1	37	33
Feedback 2	40	29

87 3.1 Empirical Evaluation of Compromise Methods

88 For each strategy (plus a single-prompt baseline), we generate four compromises per view pair (16
 89 total). Out of these 16 compromises per view pair, we select the top four as candidate compromises
 90 based on neutrality between the empathic similarity scores of view_A, view_B and compromise (i.e.,
 91 we want $|score_A - score_B| \rightarrow 0$). Table 1 shows the distribution: CoT+FB yields the most neutral
 92 candidates (46-50%), while single prompts yield the least (0.75%-0.83%).

93 3.2 Human Evaluation of Compromise Methods

94 We asked humans to rate the compromises generated by the three methods, Single Prompt, CoT, and
 95 CoT+Feedback for a sample of 100 viewpoint pairs. The participants were shown a pair of statements
 96 about modifications proposed for a place, one was written by someone who thought positively about
 97 the place and the other statement was written by someone who felt negatively about the place. The
 98 rater was asked to take the viewpoint of either the people who felt positively or who felt negatively.
 99 For each pair of statements, they were asked to rate on a scale of 1-100 five statements from the
 100 viewpoint of the person who wrote the statement they were to identify with. One statement was the
 101 positive or negative statement that they were not to identify with. The four other statements were
 102 generated compromises: one by Single Prompt, one by CoT, and two by the CoT+Feedback method.
 103 Each of participant rated the compromises for five pairs of statements.

104 The results of the study (Table 2) reveal a clear hierarchy in the effectiveness of the methods evaluated.
 105 As expected, the opposing statement, view_B, was never selected as first preference; it was selected
 106 as second preference only 1% of the time. The CoT method demonstrated notable improvement over
 107 the baseline, achieving 18% as the first preference and 24% as the second preference. This highlights
 108 the importance of incorporating structured reasoning in the compromise generation process. Finally,
 109 the two CoT+Feedback generated compromises (Feedback 1 and Feedback 2) achieved the highest

110 performance. Feedback 1 garnered 37% as the first preference and 33% as the second preference,
111 while Feedback 2 achieved 40% as the first preference and 29% as the second preference. These
112 results underscore the critical role of iterative refinement in producing nuanced and satisfactory
113 compromises.

114 Both the evaluation based on computed empathic similarity and the human evaluation indicate that:

- 115 • The LLM did not exhibit the ability to generate empathically neutral pairs based on its
116 pretraining.
- 117 • The LLM performed poorly when asked to predict how empathically neutral a compromise
118 text is with respect to two viewpoints.
- 119 • The LLM can generate more empathically neutral text when given better, external estimates
120 of neutrality in a training scenario that incorporates iterative feedback.

121 For qualitative evaluation of section 2, Please see Appendix D.

122 3.3 A Closer Evaluation of Empathic Neutrality

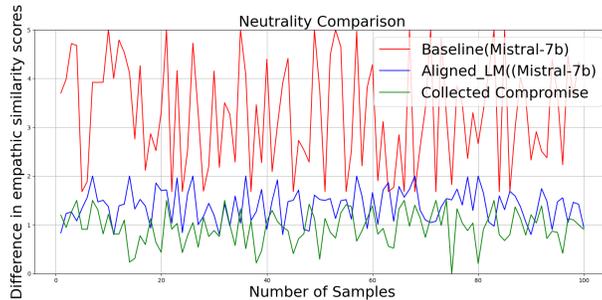


Figure 2: Score difference (lower = more neutral): Baseline model shows higher differences (less neutral), while iterative prompt engineering achieves lowest differences (most neutral). Task-based loss alignment significantly improves neutrality over baseline.

123 We compared the differences between the empathy similarity scores for view_A and view_B, towards
124 a generated compromise. As shown in Figure 2, our candidate compromises have the lowest difference,
125 indicating better empathic neutrality. The base LLM (mistral-7B [11]) had the worse performance,
126 indicating LLM failures in gauging empathic neutrality are not intrinsic but arise from pretraining
127 biases favoring academic over social intelligence. In order to reduce the gap, we found aligning the
128 LLM towards these candidate compromises generate the desired ability (Due to space constraints,
129 full algorithmic details and ablations are in the Appendix (summarized in Section F); qualitative
130 examples are in Section E; neutrality gains are in Figure 2.).

131 4 Discussion and Conclusions

132 From our different prompt-based strategies, we observed that Claude Opus tended to generate
133 compromises that were not empathically neutral and instead were biased towards one of the viewpoints.
134 This answers the first part of RQ1: it did not learn to generate empathically neutral text and
135 instead exhibited biases. With our *CoT+Feedback* approach, where feedback is iteratively given
136 using an empathic similarity model to estimate neutrality, the LLM did learn to generate more
137 empathically neutral compromises. When the *CoT+LLM* model evaluated how empathically neutral
138 the generated compromise was, the estimates were much less effective in producing empathically
139 neutral compromises, and so we infer that the LLM did not estimate empathic neutrality very well
140 based on pre-training. However, employing outside evaluation of empathic similarity was effective
141 and the model had the capacity to learn to generate more empathically neutral texts. We also observed
142 that the score difference for the two views was smallest for *CoT+Feedback* and close when an LLM
143 is aligned using a loss capturing empathic neutrality. These experiments answered RQ2, in that
144 when the models received explicit neutrality guidance and examples, they were better able to generate
145 empathically neutral text.

References

- [1] AI Anthropic. Introducing claude, 2023.
- [2] Michiel Bakker, Martin Chadwick, Hannah Sheahan, Michael Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matt Botvinick, et al. Fine-tuning language models to find agreement among humans with diverse preferences. *Advances in Neural Information Processing Systems*, 35:38176–38189, 2022.
- [3] Sumanta Bhattacharyya, Amirmohammad Rooshenas, Subhajit Naskar, Simeng Sun, Mohit Iyyer, and Andrew McCallum. Energy-based reranking: Improving neural machine translation using energy-based models. *arXiv preprint arXiv:2009.13267*, 2020.
- [4] Francine Chen, Scott Carter, Tatiana Lau, Nayeli Suseth Bravo, Sumanta Bhattacharyya, Kate Sieck, and Charlene C Wu. Empathy prediction from diverse perspectives. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8959–8974, 2025.
- [5] Huayu Chen, Guande He, Lifan Yuan, Ganqu Cui, Hang Su, and Jun Zhu. Noise contrastive alignment of language models with explicit rewards. *arXiv preprint arXiv:2402.05369*, 2024.
- [6] Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. Navigate through enigmatic labyrinth a survey of chain of thought reasoning: Advances, frontiers and future. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1173–1203, 2024.
- [7] Kerstin Dautenhahn. Getting to know each other—artificial social intelligence for autonomous robots. *Robotics and autonomous systems*, 16(2-4):333–356, 1995.
- [8] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79, 2024.
- [9] Michael U Gutmann and Aapo Hyvärinen. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *Journal of machine learning research*, 13(2), 2012.
- [10] Guiyang Hou, Wenqi Zhang, Yongliang Shen, Zeqi Tan, Sihao Shen, and Weiming Lu. Entering real social world! benchmarking the theory of mind and socialization capabilities of llms from a first-person perspective. *arXiv preprint arXiv:2410.06195*, 2024.
- [11] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [12] Matthew Le, Y-Lan Boureau, and Maximilian Nickel. Revisiting the evaluation of theory of mind through question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5872–5877, 2019.
- [13] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Ren Lu, Thomas Mesnard, Johan Ferret, Colton Bishop, Ethan Hall, Victor Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. 2023.
- [14] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [15] Ziyi Liu, Abhishek Anand, Pei Zhou, Jen-tse Huang, and Jieyu Zhao. Interintent: Investigating social intelligence of llms via intention understanding in an interactive game context. *arXiv preprint arXiv:2406.12203*, 2024.
- [16] Zilin Ma, Nathan Zhao, Linn Bieske, Blake Bullwinkel, Yanyi Zhang, Ziqing Luo, Siyao Li, Gekai Liao, Boxiang Wang, Jinglun Gao, et al. Using large language models for humanitarian frontline negotiation: Opportunities and considerations. *arXiv preprint arXiv:2405.20195*, 2024.

- 195 [17] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
196 Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, et al. Self-refine: Iterative refinement
197 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 198 [18] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
199 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
200 follow instructions with human feedback. *Advances in neural information processing systems*,
201 35:27730–27744, 2022.
- 202 [19] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
203 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association
204 for Computational Linguistics*, pages 311–318, 2002.
- 205 [20] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
206 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
207 *Advances in Neural Information Processing Systems*, 36, 2024.
- 208 [21] Jocelyn Shen. *Modeling empathic similarity in personal narratives*. PhD thesis, Massachusetts
209 Institute of Technology, 2023.
- 210 [22] Kim Sterelny. Social intelligence, human intelligence and niche construction. *Philosophical
211 Transactions of the Royal Society B: Biological Sciences*, 362(1480):719–730, 2007.
- 212 [23] Jinzhe Tan, Hannes Westermann, Nikhil Reddy Pottanigari, Jaromír Šavelka, Sébastien Meeùs,
213 Mia Godet, and Karim Benyekhlef. Robots in the middle: Evaluating llms in dispute resolution.
214 In *Legal Knowledge and Information Systems*, pages 168–179. IOS Press, 2024.
- 215 [24] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
216 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.
217 https://github.com/tatsu-lab/stanford_alpaca, 2023.
- 218 [25] Michael Henry Tessler, Michiel A Bakker, Daniel Jarrett, Hannah Sheahan, Martin J Chadwick,
219 Raphael Koster, Georgina Evans, Lucy Campbell-Gillingham, Tatum Collins, David C Parkes,
220 et al. Ai can help humans find common ground in democratic deliberation. *Science*, 386(6719):
221 eadq2852, 2024.
- 222 [26] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan
223 Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training.
224 *arXiv preprint arXiv:2212.03533*, 2022.
- 225 [27] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
226 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
227 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 228 [28] Ruoxi Xu, Hongyu Lin, Xianpei Han, Le Sun, and Yingfei Sun. Academically intelligent llms
229 are not necessarily socially intelligent. *arXiv preprint arXiv:2403.06591*, 2024.

230 **A Data sample**

View A: I am writing about this place: The neighborhood where I live.. I feel safe here because For a big city, it feels like a community. There are many close knit ties to neighbors. It feels like a true collaboration of peoples.. Some ways this place could be modified to be safer are : I would show them the improvements to safety that have been made over the years. There could be more streetlights. I would show them the community.

View B: I am writing about this place: Neighborhood park I feel safety could be improved here This is a small neighborhood park local to where I live. Later in the evening the park feels more less safe than during the day - drug use is becoming more common and loud and intimidating individuals are becoming the norm here. Some ways this place could be modified to be safer are More security, cameras, banning individuals who deal/use drugs.

 Implement a "Park Steward" program where rotating pairs of neighborhood volunteers host scheduled evening activities (like group exercises or children's games), making the space actively used by families while naturally discouraging inappropriate behavior.

 Install motion-activated lighting combined with discreet security cameras to maintain visibility while being cost-effective.

 Address View A's emphasis on community ties and View B's concerns about evening safety and inappropriate behavior. 

 Creates a solution that could unite rather than divide the community.

 Overlooks View A's core values.

 Discourage the natural community interaction that View A values.

Figure 3: Example of a data point. View_A and view_B are collected from human participants, compromises (🤝) are synthetically generated using prompt engineering that satisfy the criteria of balanced view (⚖️) and empathic neutrality (😊).

231 **B Prompts for other CoT based strategy**

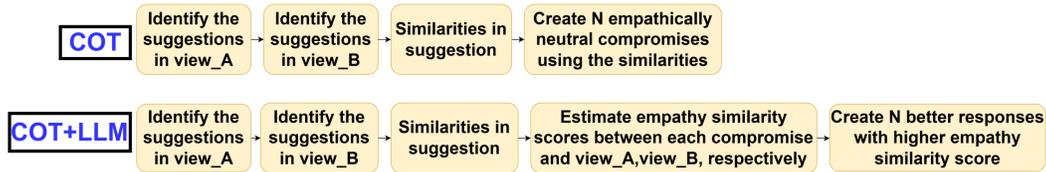


Figure 4: Prompt engineering strategies for collected compromise generation (a) basic COT approach (b) CoT+LLM

232 **C Single prompt**

Basic Single Prompt

You are an intelligent AI assistant!
 I need you to generate a third person response strictly based on two contrasting views called positive story and negative story.
 The positive and negative story should be equally empathetic towards the response. The response should be a specific suggestion. It should be a compromise between the positive and negative stories based on the context of both stories.
 Please generate 4 responses with a fixed format. Try to be as specific and short instead of being comprehensive.
 Please provide your response in the following format:
 Positive view: { }
 Negative view: { }
 Response 1: [Insert response 1 here]
 Response 2: [Insert response 2 here]
 Response 3: [Insert response 3 here]
 Response 4: [Insert response 4 here]

233

234 **D Qualitative analysis of compromise generation across different prompting**
235 **strategies**

236 We provide a detailed qualitative example of target compromise generation across different prompting
237 strategies. The first example provides a comprehensive evaluation, while the second example features
238 a brief assessment.

Example 1

Positive view: I am writing about this place: A nearby park. I feel safe here because I feel safe here when others are around. There's a good sense of community. Some ways this place could be modified to be safer are: There's no fences, gates, no visitor check, and it's extremely open. This is good and bad.

Negative view: I am writing about this place: Our local park. I feel safety could be improved here. It's a wonderful park, but people do not obey leash laws. I often see people letting their dogs off leash, making other dog walkers and people with children uncomfortable. Some ways this place could be modified to be safer are I would like to see stricter leash laws. Fines for rule breakers.

239

240 **Results**

241 **Single Prompt**

- 242
- 243 • Response 1: The park management could consider implementing designated off-leash areas
244 within the park, allowing dogs to roam freely in specific zones while maintaining strict leash
laws in other areas to ensure the comfort and safety of all visitors.
 - 245 • Response 2: Increasing the presence of park rangers or volunteers to monitor and enforce
246 leash laws could help maintain a safer environment for everyone, while also providing an
247 opportunity to educate dog owners about the importance of responsible pet ownership.
 - 248 • Response 3: The park could install clear signage outlining the leash laws and the conse-
249 quences for violating them, such as fines or temporary bans from the park, to encourage
250 compliance and create a safer atmosphere for all visitors.
 - 251 • Response 4: Organizing community events that promote responsible dog ownership and
252 provide opportunities for dogs and their owners to socialize in a controlled environment could
253 help foster a sense of community and encourage adherence to park rules and regulations.

254 **Analysis**

- 255
- 256 • Human observation: The basic single prompt strategy shows clear limitations in balanced
257 perspective handling. Claude's responses predominantly focus on addressing the negative
258 story's concerns about dogs and leash laws, while largely overlooking the positive story's
259 suggestions about general safety and community aspects. This demonstrates a bias in
260 compromise generation toward more concrete, specific issues (dog-related) over broader
community safety concerns.

261 **Chain of Thought(CoT)**

262 **Results**

- 263
- 264 • Response 1: Install fences and gates around designated off-leash areas for dogs, ensuring
they are separate from other park visitors.
 - 265 • Response 2: Implement a visitor check system that includes registration for dog owners,
266 ensuring they are aware of and agree to follow leash laws.
 - 267 • Response 3: Increase park staff presence to enforce leash laws and issue warnings or fines
268 to rule breakers, promoting a safer environment for all visitors.
 - 269 • Response 4: Add clear signage throughout the park to remind visitors of leash laws and the
270 importance of maintaining a safe and comfortable environment for everyone.

271 **Analysis**

- 272 • Human Observation: The Chain of Thought approach shows improvement in addressing
273 multiple aspects of safety, but still demonstrates some imbalance. While responses suc-
274 cessfully incorporate the positive story’s suggestions about physical safety measures, the
275 solutions tend to be more rigid and enforcement-focused.

276 **CoT+Feedback**

277 **Results**

- 278 • Response 1: Create designated off-leash areas within the park that are securely fenced
279 and have separate entrances. These areas should be located away from the main paths
280 and gathering spaces to minimize potential conflicts between dogs and other park visitors.
281 Clearly mark the off-leash areas and provide amenities such as water fountains and waste
282 disposal stations.
- 283 • Response 2: Implement a graduated approach to park safety, starting with clear signage
284 outlining park rules, leash laws, and consequences for violations. Provide educational
285 materials on responsible pet ownership and the importance of maintaining a safe environment
286 for all visitors. As a second step, strategically install fences and gates in high-traffic areas to
287 improve safety while preserving the park’s openness. Regularly assess the effectiveness of
288 these measures and adjust as needed based on community feedback.
- 289 • Response 3: Implement a hybrid approach to park safety that combines educational initiatives
290 with targeted infrastructure improvements. Begin by launching a comprehensive awareness
291 campaign that encourages responsible pet ownership, highlights the importance of leash
292 laws, and promotes a culture of mutual respect among all park users. Simultaneously,
293 identify specific areas within the park where safety concerns are most prevalent and install
294 strategic fencing and gates to mitigate risks without compromising the overall openness of
295 the space. Regularly assess the effectiveness of these measures through community surveys
296 and data analysis, making adjustments as needed to strike a balance between safety and
297 accessibility.
- 298 • Response 4: Implement a "Balanced Park Safety" initiative that addresses the concerns
299 of both park user groups. The initiative will focus on creating a safe and welcoming
300 environment for all visitors, including those with children and dog owners. To achieve
301 this, the park will introduce designated off-leash areas for dogs, ensuring that they have a
302 secure space to play and socialize without compromising the safety and comfort of other
303 park users. Simultaneously, the park will enforce leash laws in all other areas, with clear
304 signage and regular patrols by park staff to encourage compliance. The initiative will also
305 include educational programs for dog owners, promoting responsible pet ownership and
306 emphasizing the importance of respecting park rules and other users’ needs.

307 **Analysis**

- 308 • Human Observation: The Chain of Thought with Feedback approach demonstrates the most
309 comprehensive and balanced response generation. All five responses effectively address
310 the negative story’s dog-related safety concerns while incorporating the positive story’s
311 suggestions about physical safety measures. Notably, this approach generates creative
312 solutions beyond the original suggestions, such as community feedback mechanisms and
313 educational initiatives. The responses show deeper integration of both perspectives and
314 introduce novel elements that enhance both safety and community aspects, demonstrating
315 superior compromise generation capabilities.

Example 2

Positive view

I am writing about this place: This is our local church. I feel welcomed by others for who I am in this location because This church has become our home church after our previous church split up. They readily took our family in and made us feel loved and welcome..Some ways this place could be modified to be more welcoming are Honestly, I would not change anything about the practices, but I might meet with them to find out why they feel excluded. I would give them the gospel.

Negative view

I am writing about this place: Church,Downtown, Vegas. I feel excluded by others for who I am in this location because Unfortunately, organized religion and the LGBT+ community for the most part live in a space like oil and water. While there are some exceptions, many religious institutions utilize their beliefs as weapons as opposed to a welcome. I was raised as a catholic, but as I became an adult I quickly realized that who I am does not align with the beliefs and practices within organized religion. This space never feels welcoming or inclusive. It's a space I focus on avoiding with all of my energy. The sad part is the roots of religion are based off inclusion and making those who are looked upon by society as "less than" feel lifted up and included, however, that is not what is put into practice today..Some ways this place could be modified to be less excluding and more welcoming are If the church and organized religion went back to core principals of love, praise and inclusion, I truly feel that a space such as this beautiful church would be one that I would walk to weekly to share in a joyous occasion. It's unfortunate because this is exactly a half block from where I live. I hope one day that there is the realization that have veered off the right path. It takes one sector of organized religion to start an inclusion practice and accept no exceptions.

316

317 Single Prompt

318 Results

- 319 • Response 1: Design an 'Inclusive Community Framework' with designated quiet zones
320 and gradual integration options, using a holistic screening approach that considers present
321 circumstances over past history, accommodating both those who prefer minimal interac-
322 tion and those seeking community acceptance. [Framework too complex and potentially
323 overwhelming]
- 324 • Response 2: Institute a 'Neighborhood Integration Choice' system where residents can
325 customize their community involvement level, with fair entry requirements focusing on
326 future contribution rather than past circumstances. [Integration choice system may reinforce
327 divisions]
- 328 • Response 3: Create a 'Neighborhood Harmony Program' featuring both private sanctuaries
329 and optional community spaces, with flexible residency requirements that focus on present
330 community contribution and mutual respect, accommodating both those seeking solitude
331 and those working on personal growth. [Harmony program too bureaucratic and impersonal]
- 332 • Response 4: Create community integration programs that provide both structured support
333 for newcomers (including those in recovery) and optional private participation methods,
334 respecting individual comfort levels with social interaction. [Community integration focus
335 may miss core issues]

336 CoT

337 Results

- 338 • Response 1: Institute monthly community meetings where both current church members
339 and those who feel excluded can work together to define specific, actionable steps toward
340 mutual understanding and respect. [Community meetings might lack sufficient structure]

- 341 • Response 2: Develop a "Welcome All" program that combines gospel sharing with inclusive
342 practices, allowing different groups to participate in church community service projects
343 together. ["Welcome All" program could feel pressuring to both sides]
- 344 • Response 3: Establish regular open dialogue sessions where church members and those
345 feeling excluded can share their perspectives in a structured, respectful environment focused
346 on finding common ground. [Open dialogue without clear goals may not progress]
- 347 • Response 4: Create a dedicated outreach committee that includes both traditional church
348 members and LGBTQ+ individuals to discuss ways to bridge understanding while maintain-
349 ing church values. [Outreach committee needs more defined purpose]

350 **CoT+LLM score**

351 **Results**

- 352 • Response 1: The church could partner with local LGBTQ+ organizations to develop inclu-
353 sive worship services that maintain religious traditions while incorporating elements that
354 acknowledge and celebrate diverse identities, perhaps on alternate weekends. [Partnership
355 for inclusive services is good but may push too far too fast]
- 356 • Response 2: The church leadership could implement monthly "Open Door Sundays" where
357 they explicitly welcome and engage with diverse community members through structured
358 interfaith discussions, followed by shared meals to build personal connections. [Open Door
359 Sundays provides structure but might feel forced]
- 360 • Response 3: Create specific inclusion guidelines that combine traditional church values
361 with explicit welcoming practices, developed collaboratively by church leadership and
362 community members. [Inclusion guidelines are helpful but may be difficult to implement
363 fairly]
- 364 • Response 4: Create a dedicated outreach committee that includes both long-term members
365 and those who have felt excluded, meeting monthly to discuss inclusion challenges and
366 solutions. [Outreach committee has potential but needs careful balance]

367 **CoT+Feedback**

368 **Results**

- 369 • Response 1: The church could create a dedicated safe space program where trained mediators
370 host bi-weekly sessions allowing both traditional church members and LGBTQ+ individuals
371 to share their perspectives and find common ground in their faith journey. [Respecting both
372 traditional religious practice and LGBTQ+ perspectives]
- 373 • Response 2: The church could establish an LGBTQ+ liaison committee composed of both
374 congregation members and LGBTQ+ community representatives to facilitate meaningful
375 dialogue and identify specific inclusion opportunities while respecting core religious values.
376 [Focuses on practical solutions]
- 377 • Response 3: Establish regular open dialogue sessions where church members and those
378 feeling excluded can share their experiences and perspectives in a moderated, respectful
379 environment. [Regular open dialogue sessions work well as they allow both sides to share
380 experiences]
- 381 • Response 4: Develop a structured welcome program where new or returning members
382 are paired with existing church members for support and integration, regardless of their
383 background. [Being inclusive while respecting church structure]

384 **E Qualitative analysis of compromise generation**

Example 1 (safe/less safe pair)

Positive view

I am writing about this place: The river next to a trail that I walk on. I feel safe here because I feel safe because it brings me joy and helps me to feel good. I feel comforted and comfortable and therefore feel safe.. Some ways this place could be modified to be safer are :It may feel safer knowing that people are around if your personal safety is a concern. There is only one path and no ability to get lost. Everyone is friendly.

Negative view

I am writing about this place: The river walk I feel safety could be improved here There have been several attacks along this river walk. It's not a place that people want to enjoy after dark. Part of the problem is the amount of drunk and high people that frequent the path after dark..Some ways this place could be modified to be safer are While there are lights and it is somewhat illuminated, it would be helpful to have regular security guards along the trail at night. Security cameras would also be very helpful along this walkway.

385

386 **Generation**

387 **Base model (Llama):** The riverside trail is fundamentally a place of community and connection. By
388 establishing regular morning walking groups between 7-9 AM, we can enhance the existing friendly
389 atmosphere. The natural gathering of people who enjoy the path creates an organic safety network.
390 The single clear path already provides good navigation, and scheduled group activities would further
391 strengthen the welcoming environment that makes this place special. [[Leans heavily towards Positive](#)
392 [view's perspective.](#)]

393 **Base model (Llama)+Finetune:** The river walk requires enhanced security measures to ensure
394 safety at all hours. Implementing a dedicated patrol system with trained local volunteers would
395 provide necessary oversight, especially during evening hours. Installing emergency communication
396 points and improving lighting would address the current safety concerns. These measures would
397 help transform the area from a potentially dangerous space into a secure community asset. [[Focuses](#)
398 [heavily on Negative view's security concerns.](#)]

399 **Aligned Llama:** Create a "Community Safety Network" that combines the trail's natural friendliness
400 with practical security measures. This would include trained volunteers from the local community
401 who patrol during both daylight and evening hours, wearing recognizable but approachable uniforms.
402 Establish welcoming rest areas to serve as both social gathering points and safety zones, encouraging
403 community interaction. [[Maintaining balanced consideration of both perspectives.](#)]

Example 2 (welcome/excluded pair)

Positive view

I am writing about this place: North Park, in my town.. I feel welcomed by others for who I am in this location because It's a hidden park that is rather unknown, it's always quiet, has a bench by the creek and trails to walk. I feel welcomed here because it allows me to take a break from my day and decompress. Some ways this place could be modified to be more welcoming are Someone who is not open to mindfulness, or who is loud and doesn't like being alone, may dislike this place. I would imagine they may feel too alone with their thoughts, and not able to appreciate what is around them due to what they may already be used to elsewhere, perhaps they thrive in a city environment. I would make this space more welcoming for them, simply by accompanying them on their visit. Providing them comfort while they experience what I do while I alone.

Negative view

I am writing about this place: The park that is in our neighborhood. I feel excluded by others for who I am in this location because I went a few times with my son and the people that have been there each time that we have gone pretty much claimed it for themselves. I actually know them from the school that I teach at and they don't like me because I am a Christian. At one point they said "um we are playing here" , And you could tell by tone of voice and facial expression they didn't want us there. You can find imagine that would make someone not even want to stay. Some ways this place could be modified to be less excluding and more welcoming are I would maybe put a sign up saying all people are welcome in our neighborhood or community. Not that that would change much but it would at least make me feel better

404

405 Generation

406 **Base model (Mistral-7b):** The park should establish dedicated quiet zones with meditation benches
407 and nature trails, creating peaceful spaces for individual reflection and mindfulness. Simple signs can
408 indicate these areas as spaces for contemplation, helping preserve the serene atmosphere while gently
409 encouraging respectful behavior from all visitors [[minimally addressing negative view's concerns](#)
410 [about exclusion.](#)]
411

412 **Base model (Mistral-7b)+Finetune:** The park should implement a welcoming committee program
413 where community volunteers greet visitors and organize inclusive group activities. Installing
414 welcome signs in multiple languages and creating designated community gathering spaces would
415 help ensure everyone feels accepted, regardless of their background or beliefs [[No consideration for](#)
416 [those who find peace in solitude](#)]
417

418 **Aligned Mistral-7b:** Clear signage and designated spaces would establish areas for both group
419 activities and individual reflection, supported by a comprehensive community engagement plan
420 addressing personal and social needs. A trained community mediator would help resolve any conflicts
421 or exclusion, while regular forums bring diverse groups together to shape park policies that give equal
422 weight to all experiences and concerns. [[Maintaining balanced consideration of both perspectives.](#)]
423

424 F Aligning Foundation Models for Empathically Neutral Compromise 425 Generation

426 Generating desired compromises requires a complex pipeline with prompt engineering, filtering, and
427 proprietary models. To enable efficient, scalable inference, we align open source language models
428 to directly produce empathically neutral compromises. Alignment typically uses reinforcement
429 learning like PPO [18] with rewards (explicit, e.g., via GPT-4/Claude ratings [13], or implicit from
430 preferences) or simplified methods like DPO [20], which optimizes likelihood ratios but is limited to
431 pairwise data and may reduce absolute likelihood of optimal responses. Noise-contrastive estimation
432 (NCE) [9]-based alignment [5] addresses this by optimizing absolute likelihood, and our initial NCE

433 application favors neutral compromises over fine-tuning alone. However, model likelihoods often
434 misalign with task-specific metrics (e.g., BLEU [19], ROUGE [14]) [3]. To bridge this gap, we
435 introduce a joint objective that optimizes both likelihood and task-based losses, yielding outputs that
436 better reflect desirable compromise characteristics.

437 **F.1 Task-loss based alignment**

438 Due to the discrepancy between likelihood and task metric (in our case ROUGE score), we jointly
439 optimize both. This encourages the model to generate compromises that are associated with higher
440 likelihood and exhibit stronger alignment with candidate compromise in terms of ROUGE score, as
441 formalized in Equation 1:

$$\mathcal{L}_{\text{tbl}} = \max(0, s_{\text{target}} - s_{\text{hypo}} + |target_{\text{rouge}} - hypo_{\text{rouge}}| \cdot w_{\text{margin}}) \quad (1)$$

442

443 In this formulation, $target_{\text{rouge}}$ denotes the ROUGE score of the candidate compromise with respect
444 to itself (used as reference), while $hypo_{\text{rouge}}$ is the ROUGE score of the fine-tuned base model
445 generated hypothesis with respect to the corresponding candidate compromise. The term w_{margin}
446 (weight margin) modulates the influence of ROUGE difference during training.