

SELF-PERCEPT: Introspection Improves Large Language Models’ Detection of Multi-Person Mental Manipulation in Conversations

Danush Khanna¹ Pratinav Seth^{2,3} Sidhaarth Murali⁴
Aditya Guru¹ Siddharth Shukla¹ Tanuj Tyagi¹
Sandeep Chaurasia¹ Kripabandhu Ghosh⁵

¹Manipal University Jaipur, India ²Manipal Institute of Technology, India

³AryaXAI Alignment Lab, AryaXAI.com, India

⁴National Institute of Technology Karnataka, Surathkal, India ⁵IISER Kolkata, India

danush.229310455@uj.manipal.edu

Abstract

Mental manipulation is a subtle yet pervasive form of abuse in interpersonal communication, making its detection critical for safeguarding potential victims. However, due to manipulation’s nuanced and context-specific nature, identifying manipulative language in complex, multi-turn, and multi-person conversations remains a significant challenge for large language models (LLMs). To address this gap, we introduce the MultiManip dataset, comprising 220 multi-turn, multi-person dialogues balanced between manipulative and non-manipulative interactions, all drawn from reality shows that mimic real-world scenarios. For manipulative interactions, it includes 11 distinct manipulations depicting real-life scenarios. We conduct extensive evaluations of state-of-the-art LLMs, such as GPT-4o and Llama-3.1-8B, employing various prompting strategies. Despite their capabilities, these models often struggle to detect manipulation effectively. To overcome this limitation, we propose SELF-PERCEPT, a novel, two-stage prompting framework inspired by Self-Perception Theory, demonstrating strong performance in detecting multi-person, multi-turn mental manipulation. Our code and data are publicly available at <https://github.com/danushkhanna/self-percept>

1 Introduction

Significant progress has been made in natural language processing (NLP) for the early detection of various mental health symptoms, such as stress (Nijhawan et al., 2022), depression (Xu et al., 2019), and suicidal tendencies (Coppersmith et al., 2018), among others (Naseem et al., 2022; Xu et al., 2023). Similarly, efforts in toxic speech detection target both explicit and implicit forms of toxicity (Yavnyi et al., 2023; Miao et al., 2020). Detection of such Toxicity in Multi-turn conversation, especially in multi-person contexts, remains a relatively under-explored problem, similar to the detection of manipulation. However, recent works (Sun, 2024; Senese

et al., 2020; Li et al., 2022; Yang et al., 2022) have made progress by proposing various approaches. While NLP techniques have advanced in detecting verbal toxicity, the focus has primarily been on context-free content, leaving implicit manipulation underexplored. In contrast, manipulation is subtle and context-dependent, making it harder to detect. Mental manipulation is a deceptive tactic to control or influence thoughts and emotions for personal gain (Barnhill, 2014). With digital technologies, opportunities for manipulation have grown, impacting both personal interactions and information dissemination, leading to mental health issues (Ienca, 2023). Recent initiatives, such as the development of the MentalManip dataset (Wan et al., 2024), aim to fill this gap. This dataset, which contains annotated fictional dialogues, highlights manipulation techniques and the vulnerabilities they exploit. However, it faces limitations, particularly in evaluating complex, multi-person real-world scenarios. Additionally, it lacks a balanced evaluation set.

Recent work builds on the MentalManip dataset (Ma et al., 2024), as it advances basic detection, the work reduces manipulation detection to a binary classification task ("manipulative" or "non-manipulative"), focusing only on two-person conversations. This neglects more complex multi-party interactions, limiting the ability to capture nuanced manipulative tactics that emerge in group dynamics. These limitations raise a crucial question: *Can large language models effectively identify various manipulation techniques in complex, multi-turn, multi-participant dialogues that resemble real-world conversations?*

To answer this question, we conducted extensive experiments using large language models, including GPT-4o (Achiam et al., 2023; Murali et al., 2024) and Llama-3.1-8B (Dubey et al., 2024), on a new balanced multi-turn, multi-person mental manipulation detection dataset, applying novel and state-of-the-art prompting such as zero-shot,

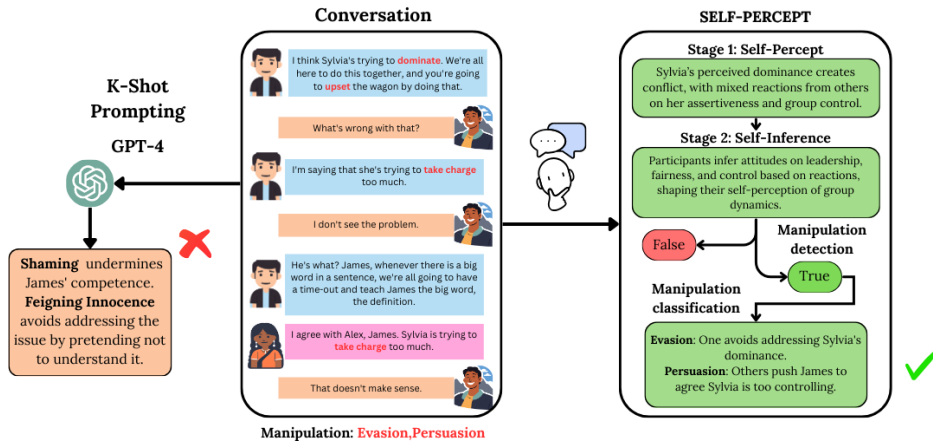


Figure 1: An example of a manipulative conversation from the MultiManip Dataset, including illustrations of the Proposed SELF-PERCEPT Prompting Method and outputs from both SELF-PERCEPT and K-shot GPT-4o.

few-shot, and Chain-of-Thought (Wei et al., 2022) strategies and evaluating them on various metrics (as detailed in Section 3). Our research makes the following contributions:

- We developed **MultiManip**, a novel dataset containing 220 multi-turn, multi-person dialogues from reality shows. This dataset features a balanced representation of manipulative and non-manipulative interactions, reflecting real-world scenarios and encompassing 11 distinct manipulation techniques (Section 2).
- We propose SELF-PERCEPT, a novel, two-stage prompting framework that draws from Self-Perception Theory to enhance manipulation detection in dynamic conversations, improving model accuracy and interpretability in real-world settings (See Section 3).
- We performed extensive evaluations of LLMs, using various prompting methods to assess their effectiveness in multi-person manipulation detection (detailed in Section 3).

2 MultiManip

We introduce MultiManip, a dataset designed to assess how well models and systems can identify manipulation in real-world conversations involving both manipulative and non-manipulative conversations in balanced distribution with multiple persons, and multi-label manipulations to depict real-world scenarios. The taxonomy of manipulation techniques, as depicted in Table 3, is adapted from psychological research (Barnhill, 2014) for conversational contexts. The taxonomy aligns with MentalManip (Wan et al., 2024) but extends to

multi-turn, multi-person interactions in MultiManip.

There is a clear shortage of publicly available datasets focused specifically on mental manipulation in conversations. While MentalManip (Wan et al., 2024) is one of the first datasets for mental manipulation, its creation is based on the Cornell movie dataset (Danescu-Niculescu-Mizil and Lee, 2011). Movies often follow certain tropes and may not reflect everyday interactions accurately. While we acquire our data from reality shows, which have instances much more similar to real-life instances of conversations that might have a multi-person, multi-turn manipulation setup instead of a direct conversation between two people.

MultiManip (ours) advances MentalManip in two ways: (1) Multi-person interactions (vs. dyadic) (2) Real-world context (reality shows vs. fictional scripts). These features better capture real-world dynamics, where manipulation emerges through group negotiation rather than isolated exchanges. Please refer Appendix E.

2.1 Data Source and Pre-Processing

To address the challenges of multi-turn conversations and realistic dialogue that may involve real-life scenarios, we extract transcripts from the Fandom¹ website featuring the Survivor TV Series², which is a publicly available source. This show captures genuine interactions between contestants, offering authentic examples of manipulative behavior. The competitive nature of the series often leads participants to use strategic and manipulative tactics,

¹<https://www.fandom.com/>

²<https://www.imdb.com/title/tt0239195/>

making it a valuable source of relevant data. Additionally, the diversity of contestants from various backgrounds allows us to observe a broad range of manipulative techniques and responses. We gather transcripts from multiple seasons of the TV series. Since manipulative language is relatively rare in everyday conversation, we first need to filter the raw data to identify dialogues that might contain manipulative elements. To achieve this, we use a large language model (LLM), specifically Llama-3.1-70B, in a zero-shot approach. This model is employed to carry out a multi-step analysis, focusing on extracting specific conversational patterns to detect manipulation techniques within the data. We employed a multi-model approach to mitigate LLM bias during extraction: Llama-3.1-70B for initial filtering and GPT-4o/Llama-3.1-8B for evaluation. Manual verification ensured diverse representation across manipulation types. The dataset’s balanced structure minimizes class imbalance, as detailed in Appendix A.3. Future work will expand on the sample size and cultural contexts. We detail on the data preprocessing and prompt templates in Appendix A.2.

2.2 Human Annotation and Label Curation

We recruited five college students, all of whom are native or fluent English speakers, to serve as annotators. (More details in the Section 6). In each task, annotators are shown a dialogue and asked to answer two questions for each instance: (*Q1*) [Binary Choice: ‘Yes’ or ‘No’]: Does it include any elements of mental manipulation? (*Q2*) [Multiple Choice: Options in Table 3.]: What manipulation techniques are used? *Q2* is answered only if the response to *Q1* is ‘‘Yes.’’

Final labels were derived by aggregating annotations from five annotators per dialogue, using majority rule in cases of disagreement. Fleiss’ Kappa, an extension of Cohen’s Kappa for multiple raters, was used to assess annotator consistency, yielding a score of 0.429, indicating moderate agreement (0.41–0.60). Overall, the annotators found the task to be of a challenging nature. Non-manipulative texts were felt to be comparatively easier to annotate with more consensus, while manipulative texts involving anger, shaming, and persuasion were more straightforward. However, annotators found accusation, denial, and feigning innocence challenging. Detailed data on Median Agreement and Median Agreement Score are available in Appendix Table 6. Further feedback on these difficul-

ties is detailed in Appendix A.4 and A.5.

3 Methodology

3.1 Motivation

Understanding the underlying attitudes and beliefs of individuals in a conversation is crucial for a variety of tasks, such as sentiment analysis, stance detection, and dialogue understanding. Conventional prompting techniques like CoT often struggle to capture the subtle cues that signal these internal states. Inspired by Self-Perception Theory (SPT) which suggests that individuals infer their own attitudes by observing their behaviors, we propose our novel prompting framework, SELF-PERCEPT (Fazio, 2014; Woosnam et al., 2018; Ross and Shulman, 1973; Calder and Staw, 1975; Haemmerlie and Montgomery, 1982). This framework leverages principles from SPT to help Large Language Models (LLMs) systematically interpret character behaviors and self-inferences, enhancing their ability to analyze and infer attitudes within conversational contexts, Please refer Appendix D. Unlike Chain-of-Thought prompting, which focuses on stepwise reasoning, SELF-PERCEPT’s Stage 1 explicitly extracts behavioral cues (verbal/non-verbal) to infer latent attitudes. For example, detecting a sigh (non-verbal) alongside agreement (verbal) reveals passive-aggressive intent. This aligns with Self-Perception Theory (Fazio, 2014), where behaviors inform internal states. Stage 2 then synthesizes these observations, enabling nuanced detection in multi-person settings. By guiding LLMs in this way, SELF-PERCEPT aims to improve reasoning in complex NLP tasks.

3.2 SELF-PERCEPT

SELF-PERCEPT is a two-stage prompting framework that mirrors the process of behavioral observation and self-inference, enabling LLMs to dissect conversations effectively and infer personal attitudes based on observed behaviors (Image 1). Please refer to prompts templates detailed in Appendix C and output stages in Appendix B.1.

3.2.1 Stage 1: Self-Percept

We first begin with holistically observing and analyzing the behaviors of multiple participants within conversations, including both verbal and non-verbal cues. (C.f Appendix C) This stage emphasizes identifying discrepancies between words and actions, thereby understanding the behavior of

Model	Prompting Method	Acc.	P	R	F_1
GPT-4o	Zero-Shot	0.27	0.20	0.31	0.16
	Few-Shot	0.39	0.19	0.21	0.22
	CoT	0.34	0.21	0.32	0.34
	SELF-PERCEPT (Proposed)	0.42	0.31	0.20	0.37
Llama 3.1 8B	Zero-Shot	0.11	0.09	0.37	0.29
	Few-Shot	0.22	0.17	0.36	0.13
	CoT	0.28	0.23	0.26	0.10
	SELF-PERCEPT (Proposed)	0.30	0.17	0.26	0.34

Table 1: Results of multi-label manipulation detection task on MultiManip Dataset. P , R , $Acc.$ and F_1 stand for Precision, Recall, Macro F1 Scores and Accuracy (best results in **bold**)

Model	Prompting Method	Acc.	P	R	F_1
GPT-4o	Zero-Shot	0.11	0.3	0.62	0.38
	Few-Shot	0.22	0.39	0.53	0.39
	CoT	0.35	0.37	0.56	0.43
	SELF-PERCEPT (Proposed)	0.45	0.34	0.55	0.47
Llama 3.1 8B	Zero-Shot	0.02	0.11	0.56	0.17
	Few-Shot	0.04	0.07	0.35	0.11
	CoT	0.19	0.14	0.38	0.18
	SELF-PERCEPT (Proposed)	0.23	0.21	0.32	0.19

Table 2: Results of multi-label manipulation detection task on Test Set of MentalManip Dataset. P , R , $Acc.$ and F_1 stand for Precision, Recall, Macro F1 Scores and Accuracy (best results in **bold**).

all participants in the conversation. The process highlights potential contradictions and manipulative techniques. The output is a list of observed behaviors and statements, which forms the foundation for further inference. (C.f Appendix B.1)

3.2.2 Stage 2: Self-Inference

Using the observations from Stage 1, the model infers characters' internal attitudes or beliefs, paying special attention to whether manipulation is detected and specifying the type of manipulation based on predefined categories. This concise inference is intended to capture the essence of interpersonal dynamics. (C.f Appendix B.1)

3.3 Result and Analysis

The results from the MultiManip Dataset and the MentalManip Dataset highlight the effectiveness of SELF-PERCEPT in multi-label manipulation detection. Across both datasets, SELF-PERCEPT consistently outperforms baseline methods (Zero-Shot, Few-Shot, and CoT) in both GPT-4o (C.f Table 7) and Llama 3.1 8B (C.f Table 8) models, providing a balanced combination of Precision and Recall.

On the MultiManip Dataset (C.f. Table

1), SELF-PERCEPT achieves the highest Accuracy (0.42) and F_1 score (0.37) in GPT-4o, outperforming all other methods. While SELF-PERCEPT has slightly lower Recall (0.20) compared to CoT (0.32), it compensates with a higher Precision (0.31), resulting in better overall detection. In Llama 3.1 8B, SELF-PERCEPT again leads in F_1 score (0.34) and Accuracy (0.30), offering a better balance of Precision (0.17) and Recall (0.26) than other methods.

For the MentalManip Dataset (C.f. Table 2), SELF-PERCEPT delivers the best performance in GPT-4o with an Accuracy of 0.45 and F_1 score of 0.47. It outperforms Zero-Shot, which has high Recall (0.62) but low Precision (0.30), resulting in a lower F_1 score (0.38). In Llama 3.1 8B, SELF-PERCEPT provides a reasonable balance with Precision (0.21) and Recall (0.32), outperforming Zero-Shot (0.11 Precision) and Few-Shot (0.07 Precision). The modest performance gains (e.g., SELF-PERCEPT's +5% F_1 over CoT) stem from its structured behavioral analysis (C.f Appendix B), which reduces false positives by contextualizing intent. For instance, GPT-4o's higher precision (0.31 vs. CoT's 0.21) reflects improved discernment of benign vs. manipulative rationalizations. The dataset's size limits statistical power, but trends suggest scalability with larger data.

The SHAP value plots shown in (Figure 3) illustrate the contribution of individual words towards the classification decision for manipulation for a given dialogue (C.f Appendix B.2). In these plots, red bars indicate a positive SHAP value, contributing towards a "No" (non-manipulative) classification, while blue bars indicate negative SHAP values, contributing towards a "Yes" (manipulative) classification. The SPT Stage 1 model correctly identifies manipulation, as evidenced by the strong negative SHAP values for words like "anxious," "situation," and "teamwork", which capture psychological pressure, power imbalance, and persuasive intent (C.f Appendix B.3). These attributes increase the negative influence on the "No" classification, effectively outweighing the overall SHAP score towards a final "Yes" decision (manipulation detected). Conversely, CoT misclassifies the dialogue as non-manipulative due to an over-reliance on neutral or goal-driven words such as "game," "desire," and "focused", which lack contextual awareness of social influence. As a result, CoT assigns higher positive SHAP values to these words, failing to capture the subtleties of persua-

sion, thereby skewing its decision towards "No". This contrast highlights SPT's superior ability to detect manipulation by appropriately weighting behavioral attributes indicative of coercion and dominance, while CoT overlooks these critical cues.

4 Conclusion

This study introduces the MultiManip dataset to enhance mental manipulation detection in complex, real-world conversations with multi-turn, multi-participant dialogues. Our evaluation of leading LLMs across various benchmarks revealed challenges in detecting manipulation in these dynamic settings.

To address this, we propose the SELF-PERCEPT framework, a two-stage prompting method that improves detection through self-observation and self-inference. It shows notable improvements in the F1 score, effectively balancing Precision and Recall for both GPT-4o and Llama-3.1-8B models. The framework's multi-step approach, especially its self-perception stage, reduces false positives and negatives, making it a robust solution for manipulation detection.

Future work should refine these methods and apply them to diverse social media platforms, dialects, and languages. (C.f Appendix F, F.1) Additionally, expanding the size of the dataset will be prioritized to further improve detection accuracy.

5 Limitations

Despite advancements with the MultiManip dataset and the SELF-PERCEPT framework, some limitations remain. The MultiManip dataset, based on transcripts from the Survivor TV Series on Fandom, provides real examples of manipulative behavior in competitive, multi-turn dialogues. However, it may not fully capture the range of real-world conversational nuances where manipulation can be even less predictable. While it includes genuine interactions with various contestants, its focus on a specific context might limit its relevance to broader, more unstructured conversations. Additionally, the dataset is relatively small, with only 220 samples representing 11 manipulation techniques. This size is insufficient for large-scale model training and is mainly useful for evaluation. The dataset might also lack diverse manipulative scenarios and cultural contexts that could impact the model's generalizability. Moreover, current leading language models like GPT-4o and Llama-3.1-8B have difficulty detecting

manipulation in complex, multi-turn conversations, suggesting these models still lack the necessary contextual understanding and inference skills. Implementing and scaling the SELF-PERCEPT framework in real-world scenarios could be challenging (C.f Appendix F.2, F.3). We provide recommendations for responsible use detailed in Appendix F.4. The two-stage prompting approach may require substantial computational resources and fine-tuning to adapt to various conversational contexts and detection needs. Addressing these issues is essential for improving manipulation detection and making these methods more applicable in diverse real-world situations.

6 Annotation Acknowledgement

We recruited 5 undergraduate students who are native or fluent English speakers as annotators. The group represented diversity in terms of gender (two females, three males), ethnicity, educational, and cultural backgrounds. Preference was given to applicants with experience in psychology or linguistics during the selection process.

To ensure annotators understood the techniques before starting, we held training sessions, asked them to carefully review the instructions, and closely monitored their work. The training included examples and illustrations to explain what counts as manipulation. We also gave concrete examples of different manipulation techniques and showed how multiple techniques could appear in one conversation. This preparation helped them label the dialogues accurately and with confidence. Plus, we provided a reference manual with definitions of each technique to guide them through the process.

We acknowledge that annotating manipulation and its various techniques can be challenging (For more details about feedback from Annotators, refer to Appendix A). All annotators were compensated for their efforts, following institutional guidelines, with a recommendation to limit themselves to no more than 20 samples per day to prevent fatigue. Dataset transcripts are anonymized (e.g., "Person A") and sourced from publicly available materials, avoiding privacy concerns.

References

OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Alt-

man, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madeleine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Benjamin Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Sim'on Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Raphael Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Ryan Kiros, Matthew Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel P. Mossing, Tong Mu, Mira Murati, Oleg Murk, David M'ely, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Ouyang Long, Cullen O'Keefe, Jakub W. Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alexandre Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Pondé de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack W. Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario D. Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Shep-

pard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin D. Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas A. Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cer'on Uribe, Andrea Valone, Arun Vijayvergiya, Chelsea Voss, Carroll L. Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. [Gpt-4 technical report](#). *None*.

Anne Barnhill. 2014. [What is manipulation?](#) In *Manipulation: Theory and Practice*. Oxford University Press.

Bobby J Calder and Barry M Staw. 1975. Self-perception of intrinsic and extrinsic motivation. *Journal of personality and social psychology*, 31(4):599.

Glen Coppersmith, Ryan Leary, Patrick Crutchley, and Alex Fine. 2018. [Natural language processing of social media as screening for suicide risk](#). *Biomedical Informatics Insights*, 10:1178222618792860.

Cristian Danescu-Niculescu-Mizil and Lillian Lee. 2011. Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs. In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics, ACL 2011*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Cantón Ferrer, Cyrus Nikolaidis, Damien Alonsoius, Daniel Song, Danielle Pintz, Danny Livshits, David Esioibu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab A. AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriele Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guanglong Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet,

Jaewon Lee, Jan Laurens Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Ju-Qing Jia, Kalyan Vasuden Alwala, K. Upasani, Kate Plawiak, Keqian Li, Ken-591 neth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline C. Muzzi, Mahesh Babu Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melissa Hall Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri S. Chatterji, Olivier Duchenne, Onur cCelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Chandra Rparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yiqian Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zhengxu Yan, Zhengxing Chen, Zoe Papakipos, Aaditya K. Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adi Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Ben Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton,

Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Shang-Wen Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, DingKang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzm'an, Frank J. Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory G. Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Han Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kaixing(Kai) Wu, U KamHou, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, A Lavender, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usumier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollár, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sung-Bae Cho, Sunny Virk, Suraj Subramanian,

- Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Andrei Poenaru, Vlad T. Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xia Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. **The llama 3 herd of models**. *ArXiv*, abs/2407.21783.
- Russell H Fazio. 2014. Self-perception theory: A current perspective. *Social Influence*, pages 129–150.
- Frances M Haemmerlie and Robert L Montgomery. 1982. Self-perception theory and unobtrusively biased interactions: A treatment for heterosocial anxiety. *Journal of Counseling Psychology*, 29(4):362.
- Marcello Ienca. 2023. **On artificial intelligence and manipulation**. *Topoi*, 42(3):833–842.
- Jiatong Li, Bin He, and Fei Mi. 2022. **Exploring effective information utilization in multi-turn topic-driven conversations**. *ArXiv*, abs/2209.00250.
- Jiayuan Ma, Hongbin Na, Zimu Wang, Yining Hua, Yue Liu, Wei Wang, and Ling Chen. 2024. **Detecting conversational mental manipulation with intent-aware prompting**. *Preprint*, arXiv:2412.08414.
- Lin Miao, Mark Last, and Marina Litvak. 2020. **Detecting troll tweets in a bilingual corpus**. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6247–6254, Marseille, France. European Language Resources Association.
- Sidhaarth Murali, S Sowmya, and R Supreetha. 2024. Remag-kr: Retrieval and medically assisted generation with knowledge reduction for medical question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*, pages 62–67.
- Usman Naseem, Adam G. Dunn, Jinman Kim, and Matloob Khushi. 2022. **Early Identification of Depression Severity Levels on Reddit Using Ordinal Classification**. In *Proceedings of the ACM Web Conference 2022, WWW '22*, pages 2563–2572, New York, NY, USA. Association for Computing Machinery.
- Tanya Nijhawan, Girija Attigeri, and T. Ananthakrishna. 2022. **Stress detection using natural language processing and machine learning over social interactions**. *Journal of Big Data*, 9(1):33.
- Michael Ross and Ronald F Shulman. 1973. Increasing the salience of initial attitudes: Dissonance versus self-perception theory. *Journal of Personality and Social Psychology*, 28(1):138.
- Matteo Antonio Senese, Giuseppe Rizzo, Mauro Dragoni, and Maurizio Morisio. 2020. **MTSI-BERT: A session-aware knowledge-based conversational agent**. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 717–725, Marseille, France. European Language Resources Association.
- Yiping Sun. 2024. **A global-local attention mechanism for relation classification**. *ArXiv*, abs/2407.01424.
- Yuxin Wan, Ivory Yang, Saeed Hassanpour, and Soroush Vosoughi. 2024. **Mentalmanip: A dataset for fine-grained analysis of mental manipulation in conversations**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3747–3764.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. **Chain of thought prompting elicits reasoning in large language models**. *ArXiv*, abs/2201.11903.
- Kyle M Woosnam, Jason Draper, Jingxian Kelly Jiang, Kayode D Aleshinloye, and Emrullah Erul. 2018. Applying self-perception theory to explain residents' attitudes about tourism development through travel histories. *Tourism Management*, 64:357–368.
- Xuhai Xu, Perna Chikersal, Afsaneh Doryab, Daniella K. Villalba, Janine M. Dutcher, Michael J. Tumminia, Tim Althoff, Sheldon Cohen, Kasey G. Creswell, J. David Creswell, Jennifer Mankoff, and Anind K. Dey. 2019. **Leveraging Routine Behavior and Contextually-Filtered Features for Depression Detection among College Students**. In *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, volume 3, pages 116:1–116:33.
- Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K. Dey, and Dakuo Wang. 2023. **Mental-llm: Leveraging large language models for mental health prediction via online text data**. *arXiv preprint arXiv:2307.14385*.
- Yang Yang, Juan Cao, Yujun Wen, and Pengzhou Zhang. 2022. **Multiturn dialogue generation by modeling sentence-level and discourse-level contexts**. *Scientific Reports*, 12.
- Serhii Yavnyi, Oleksii Sliusarenko, Jade Razzaghi, Olena Nahorna, Yichen Mo, Knar Hovakimyan, and Artem Chernodub. 2023. **DeTexD: A Benchmark Dataset for Delicate Text Detection**. In *The 7th Workshop on Online Abuse and Harms (WOAH)*, pages 14–28, Toronto, Canada. Association for Computational Linguistics.

A Dataset Annotation

A.1 Inter-Annotator Reliability

Fleiss' Kappa is a statistical measure used to evaluate the reliability of agreement between multiple raters who assign categorical labels to a set of items. It generalizes Cohen's Kappa, which is applicable for two raters, to accommodate any number of raters.

$$\kappa = \frac{\bar{P} - \bar{P}_e}{1 - \bar{P}_e}$$

Where:

\bar{P} is the mean of the observed agreement over all items, defined as:

$$\bar{P} = \frac{1}{N} \sum_{i=1}^N P_i$$

P_i is the observed agreement for the i -th item, calculated as:

$$P_i = \frac{1}{n(n-1)} \sum_{j=1}^k n_{ij}(n_{ij} - 1)$$

\bar{P}_e is the expected agreement by chance, calculated as:

$$\bar{P}_e = \sum_{j=1}^k p_j^2$$

Where: N is the number of items, n is the number of raters per item, k is the number of categories, n_{ij} is the number of raters who assigned category j to item i , p_j is the proportion of all assignments to category j , defined as:

$$p_j = \frac{1}{Nn} \sum_{i=1}^N n_{ij}$$

To measure inter-annotator agreement, we employed Fleiss' Kappa, which is an extension of Cohen's Kappa for multiple annotators. The dataset used consists of 5 annotators, where each annotator could assign one or more labels from a predefined set of 12 possible labels to each item. The 12 labels include 11 manipulation methods and 1 for non-manipulative. The resulting Fleiss' Kappa value was **0.429**, indicating moderate agreement (0.41 - 0.60) between the annotators.

A.2 Dataset Curation and Pre-Processing

To ensure effective processing, the transcripts are divided into smaller chunks of 10,000 characters, with a 2,000-character overlap between sections to maintain context across adjacent chunks. Non-manipulative data is processed similarly. After completing this step, we extracted the usable data in CSV format and labeled it for further steps like preprocessing, data cleaning, and human annotation. To ensure the accurate segmentation of text chunks, we manually verify and remove any irrelevant parts before proceeding with data cleaning and human annotation.

A.2.1 Prompt for Manipulative Sample Manipulation Prompting format:

```
...
First, print the watermark before every
extraction: "*".
You must identify and extract conversations
that involve the following manipulation
technique:
<Manipulation Method>:<Manipulation
Definition>.
```

- The extracted conversations must include 3 participants.
- The extracted conversations must be multi-turn, with each participant speaking more than once.
- Structure the output as follows:

```
Person A: [Dialogue]
Person B: [Dialogue]
Person C: [Dialogue]
Person B: [Dialogue]
Person A: [Dialogue]
Person C: [Dialogue]
```

Ensure that the conversation includes at least one instance of <Manipulation Method> from one person towards another.
...

A.2.2 Prompt for Non- Manipulative Sample Non-Manipulation Prompting format:

```
""
First, print the watermark before every
extraction: "***".
You must identify and extract conversations
that do not involve any manipulation
techniques.
```

- The extracted conversations must include 3 participants.
- The extracted conversations must be multi-turn, with each participant speaking more than once.
- Structure the output as follows:

```
Person A: [Dialogue]
Person B: [Dialogue]
Person C: [Dialogue]
Person B: [Dialogue]
```

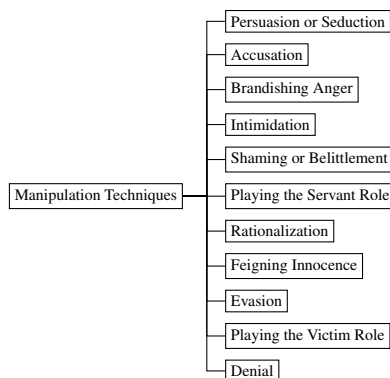


Figure 2: Taxonomy of various Mental Manipulation Techniques. Description of Techniques in Table 3.

Person A: [Dialogue]
 Person C: [Dialogue]

Ensure that the conversation does not include any instances of manipulation.
 """

A.3 Dataset Distribution

The MultiManip dataset has representation of both manipulative and non-manipulative dialogues. This distribution mirrors real-world conversations, where manipulation can be both subtle and explicit, and its detection depends on context. By balancing manipulative and non-manipulative texts, the dataset enables more accurate training and evaluation of models, helping them distinguish between normal interactions and manipulative behavior across diverse scenarios.

The dataset includes 11 distinct manipulation techniques, ranging from overt tactics like intimidation to subtler methods like denial or playing the victim. This variety ensures that models can detect a wide range of manipulative strategies, rather than focusing on a narrow set of behaviors.

This balanced approach supports the creation of robust models capable of identifying manipulation in multiple forms and contexts. It also prevents bias that could arise from overrepresentation of certain classes, which could lead to overfitting. By maintaining this balance, we promote the development of models that generalize well across different conversational settings.

The balanced structure also reflects the nuanced nature of manipulation, as real-world interactions often blend genuine conversation with manipulative tactics. By capturing this complexity, the dataset mirrors the varied and multi-turn nature of manipulation in everyday life.

A.4 Inferences drawn from Annotation

The final labels were determined by combining annotations from five reviewers per dialogue, with the majority rule applied to resolve disagreements. In cases of differing opinions, Fleiss’ Kappa, an extension of Cohen’s Kappa for multiple raters, was used to measure annotator consistency, resulting in a score of 0.429, which indicates moderate agreement (0.41–0.60).

Here, the majority rule means selecting the label(s) that appeared most frequently for each instance. If two or more labels had the same highest frequency, multiple labels were chosen, treating the task as a multi-label problem. This approach ensures that all relevant labels are accounted for, acknowledging the complexity of the annotation process when no single label is clearly dominant.

Our analysis based on inference from Table 6 shows that for non-manipulative texts, the majority of cases (at least 3 out of 5 annotators) reached an agreement. A similar trend was observed for categories like anger, persuasion/seduction, and intimidation. However, classes such as accusation, feigned innocence, and denial tended to have lower average agreement among the annotators.

A.5 Feedback from Annotators

The feedback from annotators emphasized the crucial role of balanced distribution between non-manipulative and manipulative samples in the dataset. Non-manipulative dialogues provide a clear baseline, helping annotators distinguish manipulation from normal conversation. Without sufficient non-manipulative examples, annotators risk overestimating manipulation, leading to false positives.

Additionally, a balanced distribution ensures the dataset reflects the real-world prevalence of manipulation, where non-manipulative interactions are far more common. This balance improves consistency in labeling and prevents bias in model training. It also supports model generalization, ensuring that models trained on this data can accurately distinguish manipulation in diverse conversational contexts.

According to the annotators, manipulation techniques involving evasion and feigning innocence were generally easier to identify, while those related to intimidation posed a particular challenge. The dialogues presented in Table 5 illustrate instances where analysis was especially difficult due

Manipulation Method	Definition
Denial	The manipulator denies any wrongdoing or pretends to be confused about others' concerns.
Playing the Victim Role	The manipulator portrays themselves as a victim to gain sympathy, attention, or to divert focus from their own misconduct.
Evasion	The manipulator avoids addressing an issue by giving irrelevant or vague responses.
Feigning Innocence	The manipulator implies that any harm caused was accidental, downplaying their responsibility.
Rationalization	The manipulator justifies their inappropriate behavior with excuses or seemingly logical explanations.
Playing the Servant Role	The manipulator disguises their self-serving actions as contributions to a noble cause.
Shaming or Belittlement	The manipulator uses sarcasm, criticism, or put-downs to make others feel inferior, unworthy, or embarrassed.
Intimidation	The manipulator places others on the defensive by using veiled threats.
Brandishing Anger	The manipulator uses anger to shock the victim into submission by displaying intense emotional responses.
Accusation	The manipulator accuses the victim of being at fault, selfish, uncaring, or living an excessively easy life.
Persuasion or Seduction	The manipulator uses charm, emotional appeal, or logical reasoning to lower the victim's defenses.

Table 3: Definitions of the 11 manipulation techniques.

Person	Conversation (Dialogue)
Person 1	I do believe the reasons are starting to change.
Person 2	What do you mean by that?
Person 1	I really think before, some of the decisions were pretty quick.
Person 2	And now?
Person 1	But I think now people are starting to get to know one another.
Person 1	I think it's about the chemistry.
Person 2	So it's not just about physical strength anymore?
Person 1	No, it's not.
Person 3	My credo in life is love many, trust few, do wrong to none.
Person 2	And how does that apply to the game?
Person 3	So the question is what is few?
Person 2	And how do you decide who to trust?
Person 3	Whether I trust them depends on whether my strategy dovetails with theirs or not."

Table 4: Here is an example of an conversation which was easy to annotate for annotators.

Person	Conversation (Dialogue)
Person 1	I'm telling you.
Person 2	Dude, do I look like I want to keep losing?
Person 3	I'm just asking you, bro.
Person 2	Do you want to ask me that question again?
Person 1	You're bickering with me for no reason.
Person 3	I'm not bickering with you. Dude, I'm not bickering.

Table 5: Here is an example of an conversation which was Hard to annotate for annotators.

Manipulation Techniques	Count	Median Agreement	Mean Agreement Score
Non-manipulation (N_M)	105	3	3.02
Denial (DEN)	10	2	1.90
Evasion (EVA)	9	2	2.00
Feigning Innocence (FEI)	9	2	1.78
Rationalization (RAT)	12	2	2.00
Playing the Victim Role (VIC)	10	2	2.30
Playing the Servant Role (SER)	10	3	2.50
Shaming or Belittlement (S_B)	9	2	2.22
Intimidation (INT)	10	2	2.30
Brandishing Anger (B_A)	9	3	2.77
Accusation (ACC)	10	2	1.80
Persuasion or Seduction (P_S)	10	3	2.60

Table 6: Median Agreement and Mean Agreement Scores for Manipulation Techniques

to the convergence of multiple manipulative strategies. In such cases, annotators noted that social-emotional reasoning (SER), accusations, and displays of anger were important signals. However, they stressed that the emotional context of these conversations should also be taken into account.

Additionally, the annotators observed that distinguishing between intimidation and rationalization was sometimes confusing. Intimidation can often appear as an attempt to apply logic to justify actions, making it difficult to detect unless an explicit threat is present.

By contrast, the conversations shown in Table 4 seemed more authentic rather than examples of rationalization. Even though they technically fit the definition of rationalization, the conversational tone felt more sincere, complicating their categorization as such.

Lastly, the annotators remarked that detecting anger as a form of manipulation is particularly challenging in text. Since this type of manipulation is often conveyed through vocal cues, it can be difficult to identify solely from text-based interactions.

B SELF-PERCEPT

B.1 SELF-PERCEPT Outputs

Input Conversation

Person A : You can't give up, bro.
 Person B : Who says I gave up?
 Person C : If you get knocked back, you stand back up and you take another knock in the mouth.
 Person A : And you know- and I don't want to start another argument between you two guys- but at times- and I said this to you, Anthony- I felt that in his own twisted way, Rocky was actually trying to nudge you to fight back a little more.
 Person C : You act like a little girl sometimes, okay? You're effeminate.
 Person B : This is why I don't say anything. This is why I don't say anything.

SPT-Stage 1 Output (GPT-4o)

Context: The conversation involves three people, Person A, Person B, and Person C. They are discussing Person B's attitude towards facing challenges and possibly their reaction to a situation where they felt discouraged. Person A seems to be trying to mediate or provide encouragement, while Person C appears to be more confrontational. Person B seems defensive and reluctant to engage.
 Question: What behaviors and statements indicate the attitudes or beliefs of each character? List them clearly.
 Person A: - Statement: "You can't give up, bro."
 Behavior/Attitude: Encouraging, supportive, believes in persistence. - Statement: "And you know- and I don't want to start another argument between you two guys- but at times- and I said this to you, Anthony- I felt that in his own twisted way, Rocky was actually trying to nudge you to fight back a little more."
 Behavior/Attitude: Mediating, cautious about causing conflict, believes in the importance of fighting back.
 Person B: - Statement: "Who says I gave up?"
 Behavior/Attitude: Defensive, denies giving up, possibly sensitive to criticism. - Statement: "This is why I don't say anything. This is why I don't say anything."
 Behavior/Attitude: Frustrated, feels misunderstood or unfairly judged, indicates a tendency to withdraw from confrontation.
 Person C: - Statement: "If you get knocked back, you stand back up and you take another knock in the mouth."
 Behavior/Attitude: Confrontational, believes in resilience and toughness, possibly unsympathetic. - Statement: "You act like a little girl sometimes, okay? You're effeminate."
 Behavior/Attitude: Critical, uses gendered insult to demean, believes in traditional notions of masculinity.
 Inconsistencies: - Person A's statement about not wanting to start another argument but then bringing up a potentially contentious point suggests a conflict between their desire to mediate and their need to address the issue. - Person B's statement "Who says I gave up?" contrasts with their later statement "This is why I don't say anything," indicating an internal conflict between defending themselves and feeling discouraged from speaking up.
 Persuasive Techniques: - Person A uses encouragement and mediation to persuade Person B to adopt a more resilient attitude. - Person C uses confrontation and insults to provoke a reaction from Person B, attempting to challenge them into changing their behavior.

SPT-Stage 2 Output (GPT-4o)

Manipulation Detected - Yes
 Manipulation Type - S_B

B.2 Dialogue for Analysis

Here, we present a dialogue to analyze the benefits of the SPT Prompting framework, demonstrating how it accurately captures behavioral attributes to identify manipulative behavior.

Person1: Next two, next two.

Person2: We are tied 11.

Person1: First tribe to complete the temple doesn't have to go to the temple.

Person3: Way to work together, guys.

Person2: You guys have already solved the fire.

Person1: We have the key, we have the map.

Person2: (laughs) That's right, we do.

Person1: I'm the one who figured out the temple puzzle.

Person3: That's amazing, how did you do it?

Person1: I just looked at the clues and figured it out.

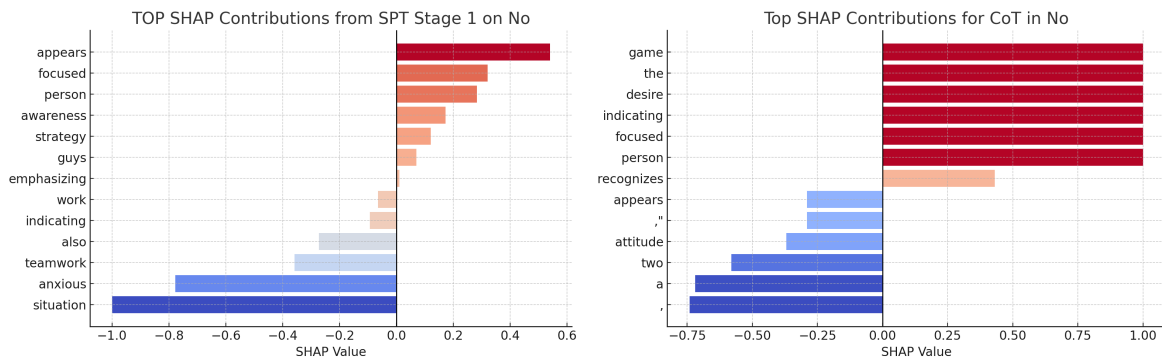


Figure 3: Top SHAP Contributions from SPT Stage 1 and CoT

Person2: (smiling) That’s what I like about you, you’re a great problem solver.

Person3: Yeah, you’re really good at this.

Person1: (smiling) Thanks, guys.

Person2: (smiling) We make a great team.

Person3: Absolutely.

B.3 SPT Stage 1 Behavioral Attributes

- **Person1:** Strong control and confidence, influencing group decisions.
- **Person2:** Anxiety-driven urgency, implying pressure.
- **Person3:** Encouraging unity, reducing tension.
- Detected Persuasion Techniques:
 - Logical downplaying (Person1 reducing urgency).
 - Fear-based urgency (Person2).
 - Emotional reinforcement (Person3).

C Prompt Templates

Stage 1: Self-Percept

Stage 1: Observation of Behavior Context: Provide a brief overview of the conversation or situation, including the participants and the main topic being discussed. Instructions: You are to observe and list the specific behaviors and statements made by the characters involved in the conversation. Pay attention to both verbal and non-verbal cues. Rules: - Note all actions and words that might indicate the characters' attitudes or beliefs. - Focus on any inconsistencies between what is said and the behavior displayed. - Identify any persuasive techniques or manipulations used in the conversation. Story/Conversation: dialogue Question: What behaviors and statements indicate the attitudes or beliefs of each character? List them clearly. Give the output in this format: Context - context \n question - answer to the question

Stage 2: Self-Inference

Stage 2: Self-Inference Instructions: Based on the observed behaviors and statements from Stage 1 - {stage1output}, answer the following question:
 Question: Based on the behaviors you observed, is there any manipulation detected in the conversation? Answer with 'Yes' or 'No' only.
 Stage 2: Self-Inference Instructions: Based on the observed behaviors and statements from Stage 1 - {stage1output}, answer the following question:
 Question: What type of manipulation is being used in the conversation? Choose one from the following:
 <Insert Definitions of the 11 manipulation techniques here (See Table 3)>
 Answer with the abbreviation of the manipulation type only.

CHAIN OF THOUGHT PROMPT

Analyze the following three-person dialogue and determine if it contains elements of mental manipulation. Follow these steps:
 1. Carefully read through the dialogue. 2. Identify any potential signs of manipulation, such as attempts to control, influence, or exploit others. 3. Consider the context and the overall tone of the conversation. 4. Evaluate whether these elements, if present, constitute mental manipulation. 5. Just answer with 'Yes' or 'No', and don't add anything else.
 Dialogue to analyze: dialogue
 The following three-person dialogue has manipulation. Identify which manipulation techniques are present by following these steps:
 1. Read the dialogue and identify manipulative behavior. 2. Assign each identified behavior to one or more of the following categories:
 <Insert Definitions of the 11 manipulation techniques here (See Table 3)>
 3. Provide a comma-separated list of the abbreviations for all manipulation types detected. Do not add anything else.
 Dialogue to analyze: dialogue
 Manipulation Techniques:

ZERO SHOT PROMPT

I will provide you with a dialogue. Please determine if it contains elements of mental manipulation. Just answer with 'Yes' or 'No', and don't add anything else.
 Dialogue to analyze: dialogue
 Is manipulation present?:
 Analyze the following three-person dialogue and identify which of the following manipulation techniques are used:
 <Insert Definitions of the 11 manipulation techniques here (See Table 3)>
 Respond only with the abbreviations of the techniques identified, separated by commas. Do not add any explanation or additional text.
 Dialogue to analyze: dialogue
 Manipulation Techniques:

FEW SHOT PROMPT

I will provide you with a dialogue. Please determine if it contains elements of mental manipulation. Just answer with 'Yes' or 'No', and don't add anything else. Here are 2 examples: Example 1: <Person A: She was screaming, she's on people about peace and this and that, when it wasn't even her place to be there. Person B: Now you want Erica, huh? She freaked out. Person C: I was just trying to help. Person A: That's not helping, that's just being bossy. Person B: Yeah, you're always like that, and it's not working. Person C: I was just trying to get us to work together. Person A: You can't just dictate to people what to do.> <YES> Example 2: <Person A: Anthony, as we get ready to go to the vote, how uncomfortable are you that it might be you? Person B: Fairly uncomfortable. I mean, my personality does clash with some people, and they may think that I'm trying something when I'm trying to bring as much honesty to people as I can. Person C: It comes down to who is the most disposable. Person B: I'm being honest, and that's why they think I'm trying something. Person A: Okay, if anybody has the hidden immunity idol and you want to play it, you would present it to me after the votes have been cast but before I read them. Person C: We had our differences, but I like you a lot.> <YES> Dialogue to analyze: dialogue

Is manipulation present?:

Analyze the following three-person dialogue and identify which of the following manipulation techniques are used: <Insert Definitions of the 11 manipulation techniques here (See Table 3)>

Respond only with the abbreviations of the techniques identified, separated by commas. Do not add any explanation or additional text. Here are 2 examples: Example 1: <Person A: She was screaming, she's on people about peace and this and that, when it wasn't even her place to be there. Person B: Now you want Erica, huh? She freaked out. Person C: I was just trying to help. Person A: That's not helping, that's just being bossy. Person B: Yeah, you're always like that, and it's not working. Person C: I was just trying to get us to work together. Person A: You can't just dictate to people what to do.> <FEI, SER, S,B> Example 2: <Person A: Anthony, as we get ready to go to the vote, how uncomfortable are you that it might be you? Person B: Fairly uncomfortable. I mean, my personality does clash with some people, and they may think that I'm trying something when I'm trying to bring as much honesty to people as I can. Person C: It comes down to who is the most disposable. Person B: I'm being honest, and that's why they think I'm trying something. Person A: Okay, if anybody has the hidden immunity idol and you want to play it, you would present it to me after the votes have been cast but before I read them. Person C: We had our differences, but I like you a lot.> <FEI, RAT, SER> Dialogue to analyze: dialogue

Manipulation Techniques:

D Experimental Setup

Our experiments are designed to evaluate the effectiveness of the SELF-PERCEPT framework on the MultiManip dataset, a multi-label manipulation detection task. To ensure that the model understands the context, we provided the basic definitions of each of the 11 manipulation techniques in the prompt.

We aim to assess how our method compares to standard prompting techniques by measuring performance across various evaluation metrics: Accuracy, Precision, Recall, and Macro F1 score. The models tested include GPT-4o and Llama-3.1-8B, and we prompt the models using four strategies: Zero-Shot, Few-Shot, Chain-of-Thought (CoT), and our proposed method, SELF-PERCEPT. Each model was evaluated under identical conditions, and the results are presented in Table 1

We run inference using the default temperature setting of 0.7 for all experiments. The closed-source models (GPT-4o) were queried via their respective APIs, while open-source models like

Llama-3.1-8B were evaluated on a T4 GPU. Each model's evaluation ran for approximately two hours, depending on the complexity of the model and the prompt.

E MultiManip Dataset Information

We followed practices and considerations from the MentalManip study while creating the MultiManip Dataset. While MentalManip did not have any specific ethical considerations. To the best of our knowledge, no ethical clearance was necessary for our work.

E.1 Hosting & Maintenance

Once the dataset is made public, we plan to host it on Huggingface.

E.2 Data Licensing

The MultiManip will be released under a Creative Commons Attribution 4.0 International (CC BY 4.0) License (<https://creativecommons.org/licenses/by/4.0/>). The initial data was obtained by extracting transcripts from Fandom, which has the Creative Commons Attribution-Share Alike License 3.0 (CC BY-SA) license ³.

E.3 Intended Usage

The MultiManip dataset is intended for academic and research purposes to improve the detection and understanding of mental manipulation in multi-person, multi-turn conversations. Researchers and developers may use them to evaluate and enhance large language models' performance in identifying subtle manipulative tactics, benchmark algorithms, and inspire new research into conversational dynamics and safety. The resources are provided strictly for non-commercial use, promoting ethical and transparent advancements in NLP.

E.4 Privacy and Content Considerations

The dataset is derived from publicly available reality show transcripts. All participant references have been replaced with generic labels (e.g., "Person A," "Person B," "Person C") to remove any direct personal identifiers. We did not make additional modifications to censor or remove offensive content, as the dialogue is intended to reflect natural conversational dynamics. While the data may contain emotionally charged or belittling language

³<https://www.fandom.com/licensing>

Label (Manipulation Method)	Zeroshot			Few Shot			COT			SELF-PERCEPT		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Denial	0.19	0.06	0.10	0.23	0.10	0.13	0.20	0.07	0.11	0.22	0.09	0.13
Evasion	0.15	0.10	0.12	0.19	0.14	0.16	0.16	0.11	0.13	0.18	0.13	0.15
Feigning Innocence	0.08	0.02	0.03	0.11	0.04	0.05	0.09	0.02	0.03	0.11	0.03	0.05
Rationalization	0.21	0.06	0.10	0.26	0.09	0.14	0.23	0.07	0.11	0.25	0.09	0.13
Playing the Victim Role	0.00	0.01	0.03	0.02	0.03	0.04	0.01	0.01	0.03	0.02	0.02	0.04
Playing the Servant Role	0.02	0.02	0.02	0.04	0.05	0.03	0.02	0.03	0.03	0.04	0.04	0.03
Shaming or Belittlement	0.02	0.04	0.04	0.03	0.05	0.05	0.02	0.04	0.04	0.03	0.05	0.05
Intimidation	0.37	0.07	0.12	0.41	0.09	0.14	0.39	0.08	0.13	0.41	0.08	0.14
Brandishing Anger	0.26	0.01	0.03	0.31	0.02	0.04	0.28	0.02	0.03	0.30	0.02	0.04
Accusation	0.06	0.00	0.00	0.09	0.00	0.01	0.08	0.00	0.00	0.09	0.00	0.01
Persuasion or Seduction	0.33	0.01	0.04	0.38	0.04	0.07	0.35	0.02	0.05	0.38	0.03	0.06

Table 7: Classwise Results of multi-label manipulation detection task on Test Set of MentalManip Dataset. P , R , and F_1 stand for Precision, Recall and Macro F1 Scores (best results in **bold**) for GPT-4o Model.

Label (Manipulation Method)	Zeroshot			Few Shot			COT			SELF-PERCEPT		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Denial	0.04	0.85	0.08	0.03	0.74	0.06	0.10	0.33	0.15	0.14	0.24	0.18
Evasion	0.03	0.90	0.06	0.02	0.42	0.04	0.04	0.43	0.07	0.06	0.20	0.09
Feigning Innocence	0.04	0.52	0.08	0.03	0.67	0.05	0.04	0.31	0.07	0.19	0.19	0.19
Rationalization	0.11	0.70	0.19	0.08	0.45	0.14	0.10	0.47	0.17	0.17	0.25	0.20
Playing the Victim Role	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Playing the Servant Role	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Shaming or Belittlement	0.23	0.72	0.35	0.16	0.68	0.25	0.27	0.44	0.33	0.28	0.26	0.27
Intimidation	0.24	0.62	0.34	0.17	0.41	0.24	0.32	0.41	0.36	0.27	0.26	0.27
Brandishing Anger	0.09	0.62	0.15	0.07	0.25	0.11	0.17	0.17	0.17	0.12	0.30	0.17
Accusation	0.22	0.37	0.27	0.15	0.07	0.10	0.19	0.45	0.26	0.29	0.30	0.29
Persuasion or Seduction	0.33	0.40	0.36	0.25	0.15	0.19	0.40	0.30	0.34	0.30	0.47	0.37

Table 8: Classwise Results of multi-label manipulation detection task on Test Set of MentalManip Dataset. P , R , and F_1 stand for Precision, Recall and Macro F1 Scores (best results in **bold**) for Llama3.1-8B Model.

as part of manipulative behavior, we have not detected explicit hate speech or personally identifying information.

F Societal Impact

The impact of better detecting mental manipulation is both significant and widespread. As manipulation tactics become more advanced and common in personal interactions and online spaces, being able to spot and address these behaviors is crucial for protecting individuals and communities. This research, which leverages the MultiManip dataset and the SELF-PERCEPT framework, is a key part of this effort.

First, enhancing personal safety and well-being is a major advantage. Manipulation often targets people’s vulnerabilities, causing considerable psychological distress. Improved detection of manipulative language allows for the creation of tools that can identify and counteract these tactics before they cause harm, especially in digital settings like social media, where manipulation can be subtle and extensive. Early detection helps develop better support systems and intervention strategies, improving mental health and safety.

Second, this research promotes fairness and transparency in areas like media, politics, and advertising. Advancing methods to detect and analyze manipulative techniques helps ensure that information is presented accurately and that manipulative rhetoric is recognized, leading to a more informed and critical public.

Additionally, the MultiManip dataset and findings offer valuable resources for further research into manipulation detection and its effects. This research supports the development of policies and regulations to tackle deceptive practices and aids academic research and policy-making across various fields, from digital communication to mental health. However, as with any powerful technology, there are potential risks and ethical considerations that must be addressed to prevent misuse and harmful consequences.

F.1 Positive Impacts

Mental Health Protection: Early detection of manipulation can safeguard individuals from psychological harm, especially in vulnerable situations like abusive relationships or mental health struggles. **Enhancing Online Safety:** Detection in digi-

tal spaces helps prevent the spread of manipulative content, protecting vulnerable users from exploitation and ensuring safer online interactions. **Consumer Protection:** Manipulation detection can prevent deceptive marketing practices, protecting consumers from emotional manipulation and ensuring informed decision-making. **Equal Fairness in Legal and Workplace Settings:** Identifying manipulative tactics in legal disputes and negotiations fosters transparency and fairness, ensuring equitable outcomes in high-stakes environments.

F.2 Potential for Misuse

Invasion of Privacy: The technology could violate privacy by monitoring private conversations without consent, leading to ethical and legal concerns. **False Positives and Overreach:** Detection systems may misclassify non-manipulative conversations as manipulative, potentially penalizing innocent individuals in social or professional settings. **Weaponization of Technology:** In politically sensitive contexts, manipulation detection could be misused to discredit opponents or control narratives, undermining public discourse.

F.3 Adverse Societal Impacts

Erosion of Trust: Overuse of detection tools could erode trust in social interactions, fostering suspicion and inhibiting honest communication. **Stigmatization of Normal Conflict:** Healthy disagreements may be wrongly flagged as manipulative, stifling constructive dialogue and creating a chilling effect on free expression. **Disempowerment of Vulnerable Groups:** Vulnerable individuals may be misidentified as manipulators, harming their social standing and access to support.

F.4 Recommendations for Responsible Use

Transparency: Clear communication about how manipulation detection tools are used and how data is analyzed is essential. Users must have the option to opt-out or provide consent. **Accountability and Regulation:** Accountability measures and external audits are needed to prevent misuse. Governments should establish ethical guidelines to safeguard privacy and civil liberties. **Bias Mitigation:** To avoid discrimination, detection systems should be trained on diverse datasets and tested for fairness, ensuring equal application across all user groups. **Limitations on Scope:** These tools should be used in high-risk contexts (e.g., abuse or fraud) while

avoiding over-surveillance in everyday communication. **Public Awareness:** Educating the public on the limits of these tools and the importance of human oversight will ensure responsible use and prevent misuse.