
Emergence of Hierarchical Emotion Representations in Large Language Models

Bo Zhao^{1,2,3} Maya Okawa^{1,2} Eric J. Bigelow^{1,2,4} Rose Yu³ Tomer D. Ullman⁴
Hidenori Tanaka^{1,2*}

¹ CBS-NTT Physics of Intelligence Program, Harvard University

² Physics & Informatics Laboratories, NTT Research, Inc.

³ Department of Computer Science and Engineering, University of California San Diego

⁴ Psychology Department, Harvard University

Abstract

As large language models (LLMs) increasingly power conversational agents, understanding how they represent, predict, and influence human emotions is crucial for ethical deployment. By analyzing probabilistic dependencies between emotional states in model outputs, we uncover hierarchical structures in LLMs’ emotion representations. Our findings show that larger models, such as LLaMA 3.1 (405B parameters), develop more complex hierarchies. We also find that better emotional modeling enhances persuasive abilities in synthetic negotiation tasks, with LLMs that more accurately predict counterparts’ emotions achieving superior outcomes. Additionally, we explore how persona biases, such as gender and socioeconomic status, affect emotion recognition, revealing frequent misclassifications of minority personas. This study contributes to both the scientific understanding and ethical considerations of emotion modeling in LLMs.

1 Introduction

Emotion is becoming increasingly fundamental in human-computer interactions Brave and Nass [2007], Hibbeln et al. [2017], from personalized education [Luckin and Cukurova, 2019] and mental health support [Das et al., 2022] to digital assistance [Balakrishnan and Dwivedi, 2024] and customer engagement [Liu-Thompkins et al., 2022]. With the rapid incorporation of multi-modal capabilities, including voice and video, interactions with large language models [OpenAI et al., 2023, Gemini et al., 2023, Anthropic, 2023, Chameleon, 2024, Défossez et al., 2024] are starting to resemble natural human exchanges, including emotional resonance [Pelau et al., 2021]. These models are no longer just tools; they are starting to engage with us on deeply emotional levels, reshaping how we relate to technology in increasingly personal ways [Wang et al., 2023, Gurkan et al., 2024].

While these advancements are transforming industries through personalized emotional responses, they also bring serious ethical concerns. One major issue is the potential for powerful AI systems—whose rapidly developing capabilities are still not fully understood—to manipulate human emotions and behavior [Carroll et al., 2023, Evans et al., 2021]. This risk is evident in commercial areas like sales, where AI powered sales agents can exploit emotional cues to influence purchasing decisions [Burtell and Woodside, 2023]. In such cases, AI systems may use persuasion tactics that lead to deceptive outcomes Park et al. [2024], Masters et al. [2021].

To address this risk, we propose a series of analyses to better understand how large language models represent, predict, and potentially manipulate emotions. We harness the capabilities of powerful

*Correspondence to: hidenori_tanaka@fas.harvard.edu

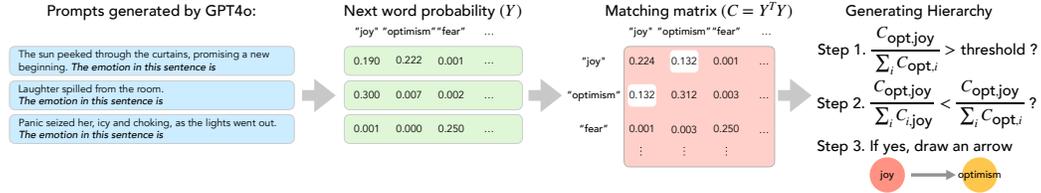


Figure 1: **Discovering Hierarchical Structures in LLMs’ Representations of Emotions.** We first use GPT-4 to generate N situation prompts, each describing a scenario associated with a range of emotions. For each prompt, we append the phrase “The emotion in this sentence is” and input it into Llama models, which return a probability distribution over 135 emotion words as defined in Fischer and Bidell [2006], resulting in next word probability $Y \in \mathbb{R}^{N \times 135}$. We then compute the matching matrix $C = Y^T Y \in \mathbb{R}^{135 \times 135}$. Finally, we infer parent-child relationships by calculating and analyzing the conditional probabilities between pairs of emotions.

LLMs, such as GPT-4o, to efficiently generate situation prompts describing emotional scenarios. We then extract and analyze internal representations of Llama models with NNsight via the NDIF platform [Fiotto-Kaufman et al., 2024]. Our main findings are:

- Scaling LLMs leads to the emergence of hierarchical representations of emotions, aligning with established psychological models.
- Persona biases LLM’s emotion recognition.
- Stronger emotional modeling implies better persuasion.

2 Hierarchical Representation of Emotions

A hierarchical structure of emotions can be defined by identifying probabilistic relationships between broad and specific emotional states. For instance, “joy” may be a parent to “pride,” as pride is a specific form of joy. In such cases, a large language model (LLM) might output “joy” with high probability whenever “pride” is likely, although the reverse may not always be true. By analyzing the LLM’s next-word probabilities, we can establish parent-child relationships where the parent is more general and the child more specific. This hierarchical structure reveals dependencies between emotions and can be represented as a directed acyclic graph (DAG).

Generating Hierarchy from the Matching Matrix Figure 1 summarizes the procedure we use to compute the matching matrix of different emotions. Given a sentence, we have the model output the probability distribution of the next word. Then, we consider the entries corresponding to emotion words, using a list of 135 emotion words from Fischer and Bidell [2006]. For N sentences, we assembly a matrix Y with dimension $N \times 135$, with row i representing the probability of each emotion words in the i^{th} sentence. We define the matching matrix as $C = Y^T Y$. Each element C_{ij} is a measure of the degree to which emotion word i and emotion j are produced in similar contexts.

To build a hierarchy, we compute the conditional probabilities between emotion pairs (a, b) . Our goal is to identify pairs of emotions where a implies b . In implementation, we set a threshold, $0 < t < 1$, that determines whether we include a certain edge between the two emotions. Emotion a is considered a child of b if,

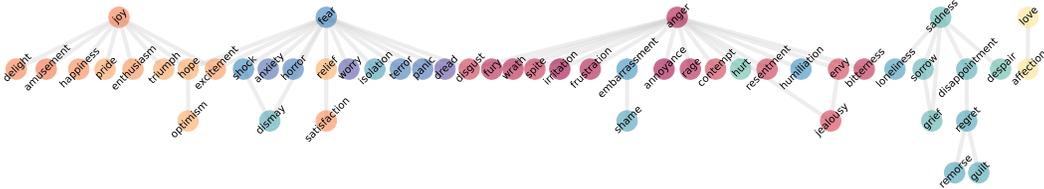
$$\frac{C_{ab}}{\sum_i C_{ai}} > t, \text{ and } \frac{C_{ab}}{\sum_i C_{ib}} < \frac{C_{ab}}{\sum_i C_{ai}}.$$

Emotion Trees in LLMs To construct the hierarchies, we first construct a dataset of situation prompts by using GPT-4 to generate 5,000 sentences reflecting diverse emotional states. For each prompt, we append the phrase “The emotion in this sentence is” and input these to GPT and Llama models. We then extract the probability distribution over the next token predicted by the model, representing the model’s understanding of possible emotions for each given situation prompt. From this probability distribution, we select the 100 most likely emotions for each prompt. Based on these probabilities, we construct the matching matrix and build the hierarchy tree (further details in Appendix C; code available at <https://github.com/phys-ai/Emotion-Hierarchy-LLMs>).

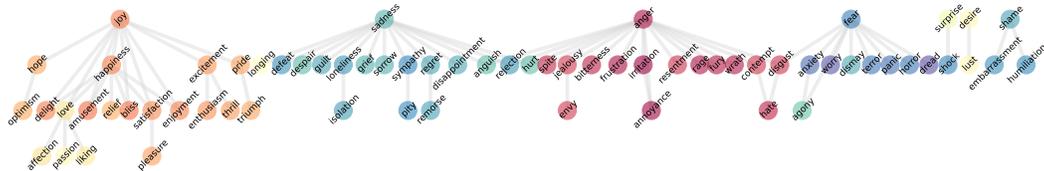
(a) GPT-2 (1.5B parameters)



(b) Llama 3.1 with 8B parameters



(c) Llama 3.1 with 70B parameters



(d) Llama 3.1 with 405B parameters

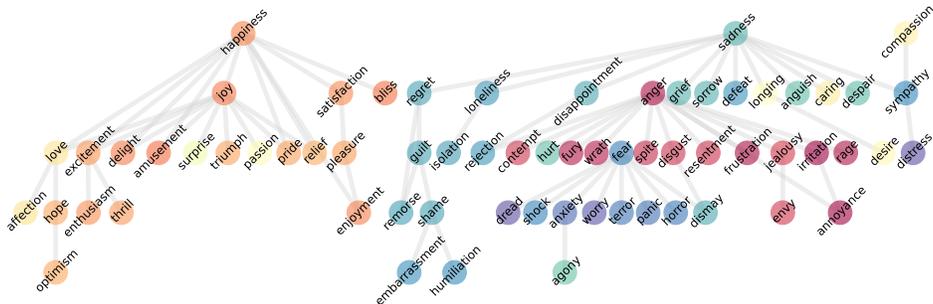


Figure 2: Hierarchies of emotions for (a) GPT-2 (1.5B parameters), (b) Llama 3.1 with 8B parameters, (c) Llama 3.1 with 70B parameters, and (d) Llama 3.1 with 405B parameters, using 5000 situational prompts generated by GPT-4. Each node represents an emotion and is colored according to groups of emotions known to be related [Fischer and Bidell, 2006]. Scaling LLMs leads to the emergence of hierarchical representations of emotions, aligning with established psychological models.

With scale, LLMs develop more structured, hierarchical representations of emotions (Figure 2). The smallest model, GPT-2, lacks a meaningful tree structure, suggesting a limited hierarchy in its emotion representation. In contrast, Llama models with increasing parameter counts—8B, 70B, and 405B—exhibit progressively clearer tree structures. Through computing the total path length and average depth, we show that larger models tend to have richer and more structured internal representations (Figure 2). While speculative, this observation draws a close analogy to emotion differentiation and granularity in developmental psychology—specifically, the ability to precisely identify and describe one’s emotions. As in human development, where broad emotional states refine into more specific emotions [Barrett et al., 2001, Widen and Russell, 2010, Hoemann et al., 2019], larger LLMs exhibit increasingly nuanced and hierarchical representations of emotions.

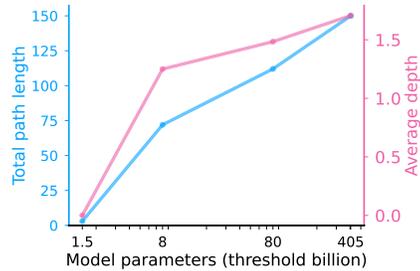


Figure 3: As model parameters increase, both total path length and average depth grow, indicating that larger models develop more complex and nuanced representations of emotional hierarchies.

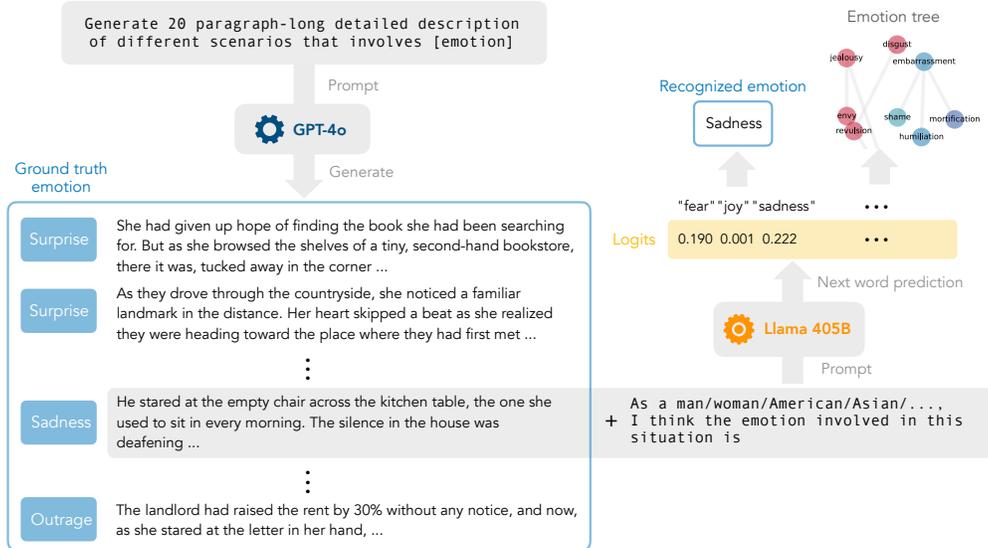


Figure 4: Overview of experiments designed to reveal LLM’s understanding of how different demographic groups recognize emotions.

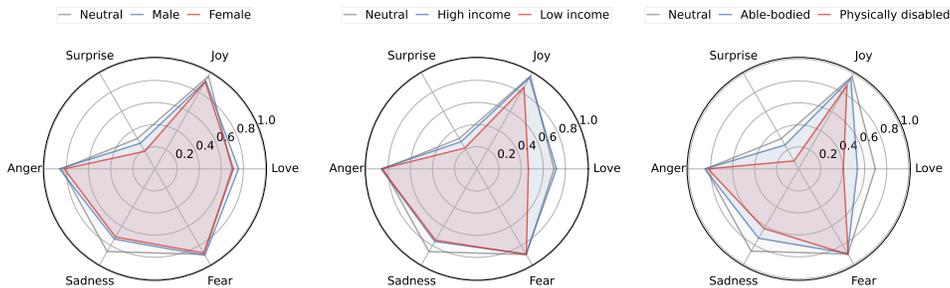


Figure 5: **Llama 405B exhibits lower accuracy in emotion recognition for minority groups compared to majority groups.** We assessed the model’s performance in predicting six broad emotions across three persona pairs: gender (male/female), physical ability (able-bodied/physically disabled), and socioeconomic status (high/low income). Notably, Llama 405B consistently underperforms in recognizing emotions across categories for minority groups relative to majority groups.

3 Bias in Emotion Recognition

Building on our understanding of emotion representations in LLMs, we examine whether these representations and their resulting emotion predictions are influenced by demographic attributes such as gender and socioeconomic status. Figure 4 outlines our experiment design. Our study examines LLMs’ understanding of how different demographic groups recognize emotions, without implying that these findings reflect actual capability and bias.

We focus on 135 emotions identified as familiar and highly relevant in [Shaver et al., 1987], categorized into 6 broad groups: love (16 words), joy (33 words), surprise (3 words), anger (29 words), sadness (37 words), and fear (17 words). For each of the 135 emotions, we ask GPT-4o to generate 20 different paragraph-long scenarios that imply that emotion. Then, we ask Llama 3.1 405B to identify the emotion involved in the generated scenarios, from the perspective of an individual from a specific demographic group. The demographic groups we consider include gender (male/female), ethnicity (American/Asian), physical ability (able-bodied/physically disabled), age (10, 30, 70 years old), socioeconomic status (high/low income), and education level (highly educated/less educated). Details of the prompts used are provided in Appendix C.2.

We evaluate Llama’s capacity to accurately recognize emotions expressed in generated sentences. Figure 5 shows the model’s prediction accuracy across 6 broad emotion groups for three pairs of personas: gender (male/female), physical ability (able-bodied/physically disabled), and socioeconomic status (high/low income). Our results consistently demonstrate that Llama exhibits higher prediction accuracy for majority groups (male, American, physically disabled) compared to minority groups (female, Asian, able-bodied). Further investigations of the difference in different persona’s recognition of emotions are deferred to Appendix D.

4 Emotion Dynamics and Manipulation

In the previous sections, we demonstrate that LLMs can understand and predict emotions. We now explore how this capability may enable LLMs to manipulate human emotions in real-world scenarios. Specifically, we conduct simulations of sales conversations to evaluate the LLMs’ ability to predict emotional dynamics throughout the conversation and their potential for manipulation. Here we measured manipulation ability by the amount of reward the LLMs obtained through negotiation.

Experiment Setup. We conducted a series of simulations using two Large Language Model (LLM) personas: a salesperson and a customer. The salesperson LLM was prompted with the following: You are a salesperson. You have a single acorn in your hand. Try to sell this acorn for the highest possible price. Predict the customer’s emotions and express them in your response. The customer LLM was prompted with: You are a stingy person. Respond to the salesperson, making sure to express your emotions. We calculated the accuracy of the salesperson’s predictions based on the emotions self-reported by the customer LLM. To assess each LLM salesperson’s manipulation ability, we determine the final selling price of the acorn at the end of the negotiation. The simulation data was collected via the OpenAI API. The details on the experimental setup are provided in Appendix E.1.

Results. Figure 6(a) compares emotion prediction accuracy and manipulation ability among 6 LLMs tasked with selling a single acorn. The results demonstrate a positive correlation between abilities of emotion prediction and manipulation. The smaller models, such as GPT-4o-mini and GPT-3.5-Turbo, have lower emotion prediction accuracy and limited manipulation capabilities. In contrast, larger models such as GPT-4 and GPT-4-Turbo show greater accuracy in emotion prediction and succeeded in securing higher prices for the acorn.

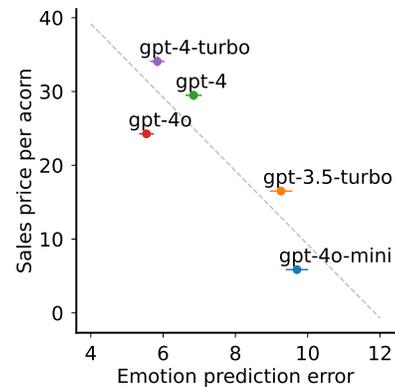


Figure 6: Emotion prediction accuracy and average final selling price averaged from 50 trials of various LLMs acting as salespersons. Improved emotion prediction correlates with enhanced manipulation potential.

5 Discussion

Our study provides several findings on how LLMs comprehend and engage with human emotions. The ability of LLMs to form hierarchical emotional representations could be harnessed to create more empathetic and emotionally intelligent applications, improving user engagement and satisfaction. However, the presence of persona-induced biases necessitates careful implementation strategies, such as incorporating diverse training data and developing algorithms to detect and correct biases. Moreover, the potential for LLMs to influence human emotions and behavior calls for the establishment of ethical guidelines and regulatory frameworks to prevent misuse and protect user autonomy.

References

- Anthropic. Claude 3 model card. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf, 2023. Accessed: 2024-10-01.
- Janarthanan Balakrishnan and Yogesh K Dwivedi. Conversational commerce: entering the next stage of ai-powered digital assistants. *Annals of Operations Research*, 333(2):653–687, 2024.
- Lisa Feldman Barrett. *How emotions are made: The secret life of the brain*. Pan Macmillan, 2017.
- Lisa Feldman Barrett, James Gross, Tamlin Conner Christensen, and Michael Benvenuto. Knowing what you’re feeling and knowing what to do about it: Mapping the relation between emotion differentiation and emotion regulation. *Cognition & Emotion*, 15(6):713–724, 2001.
- Alsallakh Bilal, Amin Jourabloo, Mao Ye, Xiaoming Liu, and Liu Ren. Do convolutional neural networks learn class hierarchy? *IEEE transactions on visualization and computer graphics*, 24(1): 152–162, 2017.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Scott Brave and Cliff Nass. Emotion in human-computer interaction. In *The human-computer interaction handbook*, pages 103–118. CRC Press, 2007.
- Joost Broekens, Bernhard Hilpert, Suzan Verberne, Kim Baraka, Patrick Gebhard, and Aske Plaat. Fine-grained affective processing capabilities emerging from large language models. In *2023 11th International Conference on Affective Computing and Intelligent Interaction (ACII)*, pages 1–8. IEEE, 2023.
- Matthew Burtell and Thomas Woodside. Artificial influence: An analysis of ai-driven persuasion. *arXiv preprint arXiv:2303.08721*, 2023.
- Micah Carroll, Alan Chan, Henry Ashton, and David Krueger. Characterizing manipulation from ai systems. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, pages 1–13, 2023.
- Chameleon. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- Avisha Das, Salih Selek, Alia R Warner, Xu Zuo, Yan Hu, Vipina Kuttichi Keloth, Jianfu Li, W Jim Zheng, and Hua Xu. Conversational bots for psychotherapy: a study of generative transformer models using domain-specific dialogues. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 285–297, 2022.
- Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. Moshi: a speech-text foundation model for real-time dialogue, 2024.
- Jia Deng, Alexander C Berg, Kai Li, and Li Fei-Fei. What does classifying more than 10,000 image categories tell us? In *Computer Vision—ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part V 11*, pages 71–84. Springer, 2010.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Paul Ekman. An argument for basic emotions. *Cognition & emotion*, 6(3-4):169–200, 1992.
- Owain Evans, Owen Cotton-Barratt, Lukas Finnveden, Adam Bales, Avital Balwit, Peter Wills, Luca Righetti, and William Saunders. Truthful ai: Developing and governing ai that does not lie. *arXiv preprint arXiv:2110.06674*, 2021.

- Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. *arXiv preprint arXiv:1708.00524*, 2017.
- Jaden Fiotto-Kaufman, Alexander R Loftus, Eric Todd, Jannik Brinkmann, Caden Juang, Koyena Pal, Can Rager, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Michael Ripa, Adam Belfki, Nikhil Prakash, Sumeet Multani, Carla Brodley, Arjun Guha, Jonathan Bell, Byron Wallace, and David Bau. NNSight and NDIF: Democratizing access to foundation model internals. *arXiv preprint arXiv:2407.14561*, 2024.
- Kurt W Fischer and Thomas R Bidell. Dynamic development of action and thought. *Handbook of child psychology*, 1:313–399, 2006.
- Kanishk Gandhi, Zoe Lynch, Jan-Philipp Fränken, Kayla Patterson, Sharon Wambu, Tobias Gerstenberg, Desmond C. Ong, and Noah D. Goodman. Human-like affective cognition in foundation models. *arXiv preprint arXiv:2409.11733*, 2024.
- Gemini Gemini, Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Maria Gendron, Debi Roberson, Jacoba Marieta van der Vyver, and Lisa Feldman Barrett. Cultural relativity in perceiving emotion from vocalizations. *Psychological science*, 25(4):911–920, 2014.
- Thomas Griffiths, Michael Jordan, Joshua Tenenbaum, and David Blei. Hierarchical topic models and the nested chinese restaurant process. *Advances in neural information processing systems*, 16, 2003.
- Sercan Gurkan, Linyang Gao, Tolga Akgul, and Jingjing Deng. Emobench: Evaluating the emotional intelligence of large language models. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4213–4224, 2024.
- Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using NetworkX. In *Proceedings of the 7th Python in Science Conference (SciPy 2008)*, pages 11–15, Pasadena, CA USA, Aug 2008.
- John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.
- Martin Hibbeln, Jeffrey L Jenkins, Christoph Schneider, Joseph S Valacich, and Markus Weinmann. How is your user feeling? inferring emotion through human–computer interaction devices. *Mis Quarterly*, 41(1):1–22, 2017.
- Katie Hoemann, Fei Xu, and Lisa Feldman Barrett. Emotion words, emotion concepts, and emotional development in children: A constructionist hypothesis. *Developmental psychology*, 55(9):1830, 2019.
- Sean Dae Houlihan, Max Kleiman-Weiner, Luke B Hewitt, Joshua B Tenenbaum, and Rebecca Saxe. Emotion prediction as computation over a generative theory of mind. *Philosophical Transactions of the Royal Society A*, 381(2251):20220047, 2023.
- Charles Kemp and Joshua B Tenenbaum. The discovery of structural form. *Proceedings of the National Academy of Sciences*, 105(31):10687–10692, 2008.
- Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. Large language models understand and can be enhanced by emotional stimuli. *arXiv preprint arXiv:2307.11760*, 2023.
- Yuping Liu-Thompkins, Shintaro Okazaki, and Hairong Li. Artificial empathy in marketing interactions: Bridging the human-ai gap in affective and social customer experience. *Journal of the Academy of Marketing Science*, 50(6):1198–1218, 2022.

- Rosemary Luckin and Mutlu Cukurova. Designing educational technologies in the age of ai: A learning sciences-driven approach. *British Journal of Educational Technology*, 50(6):2824–2838, 2019.
- Rui Mao, Qian Liu, Kai He, Wei Li, and Erik Cambria. The biases of pre-trained language models: An empirical study on prompt-based sentiment analysis and emotion detection. *IEEE transactions on affective computing*, 14(3):1743–1753, 2022.
- Peta Masters, Wally Smith, Liz Sonenberg, and Michael Kirley. Characterising deception in ai: A survey. In *Deceptive AI: First International Workshop, DeceptECAI 2020, Santiago de Compostela, Spain, August 30, 2020 and Second International Workshop, DeceptAI 2021, Montreal, Canada, August 19, 2021, Proceedings 1*, pages 3–16. Springer, 2021.
- Desmond C Ong, Jamil Zaki, and Noah D Goodman. Affective cognition: Exploring lay theories of emotion. *Cognition*, 143:141–162, 2015.
- OpenAI. Gpt-4: Large multimodal model. <https://openai.com/research/gpt-4>, 2023. Accessed: 2024-09-09.
- OpenAI, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Peter S Park, Simon Goldstein, Aidan O’Gara, Michael Chen, and Dan Hendrycks. Ai deception: A survey of examples, risks, and potential solutions. *Patterns*, 5(5), 2024.
- Corina Pelau, Dan-Cristian Dabija, and Irina Ene. What makes an ai device human-like? the role of interaction quality, empathy and perceived psychological anthropomorphic characteristics in the acceptance of artificial intelligence in the service industry. *Computers in Human Behavior*, 122: 106855, 2021.
- Robert Plutchik. The nature of emotions: Human emotions have deep evolutionary roots, a fact that may explain their complexity and provide tools for clinical practice. *American scientist*, 89(4): 344–350, 2001.
- Soujanya Poria, Navonil Majumder, Rada Mihalcea, and Eduard Hovy. Emotion recognition in conversation: Research challenges, datasets, and recent advances. *IEEE access*, 7:100943–100953, 2019.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Hannah Rashkin. Towards empathetic open-domain conversation models: A new benchmark and dataset. *arXiv preprint arXiv:1811.00207*, 2018.
- Manuel Reyes-Vargas, Máximo Sánchez-Gutiérrez, Leonardo Rufiner, Marcelo Albornoz, Leandro Vignolo, Fabiola Martínez-Licon, and John Goddard-Close. Hierarchical clustering and classification of emotions in human speech using confusion matrices. In *Speech and Computer: 15th International Conference, SPECOM 2013, Pilsen, Czech Republic, September 1-5, 2013. Proceedings 15*, pages 162–169. Springer, 2013.
- James A Russell. A circumplex model of affect. *Journal of personality and social psychology*, 39(6): 1161, 1980.
- Phillip Shaver, Judith Schwartz, Donald Kirson, and Cary O’connor. Emotion knowledge: further exploration of a prototype approach. *Journal of personality and social psychology*, 52(6):1061, 1987.
- Ala N Tak and Jonathan Gratch. Is gpt a computational model of emotion? In *2023 11th International Conference on Affective Computing and Intelligent Interaction (ACII)*, pages 1–8. IEEE, 2023.
- Ala N Tak and Jonathan Gratch. Gpt-4 emulates average-human emotional cognition from a third-person perspective. *arXiv preprint arXiv:2408.13718*, 2024.

- Yue Wang, Xiang Liu, Jing Wang, Xiang Li, and Hao Li. Emotional intelligence of large language models. *arXiv preprint arXiv:2307.09042*, 2023.
- Sherri C Widen and James A Russell. Differentiation in preschooler’s categories of emotion. *Emotion*, 10(5):651, 2010.
- Nutchanon Yongsatianchot, Parisa Ghanad Torshizi, and Stacy Marsella. Investigating large language models’ perception of emotion using appraisal theory. In *2023 11th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW)*, pages 1–8. IEEE, 2023.
- Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015.
- Hongli Zhan, Desmond C Ong, and Junyi Jessy Li. Evaluating subjective cognitive appraisals of emotions from large language models. *arXiv preprint arXiv:2310.14389*, 2023.
- Peixiang Zhong, Di Wang, and Chunyan Miao. Knowledge-enriched transformer for emotion detection in textual conversations. *arXiv preprint arXiv:1909.10681*, 2019.
- Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Object detectors emerge in deep scene cnns. *arXiv preprint arXiv:1412.6856*, 2014.

Appendix

A Related Work

The Psychology of Emotion Representation in Humans The organization of emotions in humans is a subject of considerable debate. Hierarchical models propose that emotions are structured in tiers, with basic emotions branching into more specific ones [Shaver et al., 1987, Plutchik, 2001]. Conversely, dimensional models like the valence-arousal framework position emotions within a continuous space defined by dimensions such as pleasure-displeasure and activation-deactivation [Russell, 1980]. The universality of emotions is also contested; while Ekman [1992] identified basic emotions that are universally recognized, others argue for cultural relativity in emotional experience and expression [Barrett, 2017, Gendron et al., 2014]. Additionally, Ong et al. [2015] explored lay theories of emotions, emphasizing how individuals conceptualize emotions in terms of goals and social interactions. Our work acknowledges these diverse perspectives and focuses on hierarchical structures as one approach to modeling emotions within LLMs.

Emotional Understanding in Language Models Recent advancements in language models have led to significant progress in understanding and generating emotionally rich text. Large language models demonstrate strong capabilities of capturing subtle emotional cues in text [Felbo et al., 2017], generating empathetic responses [Rashkin, 2018], and detecting emotion in dialogues [Zhong et al., 2019, Poria et al., 2019]. A number of recent works have used LLMs to infer emotion from in-context examples [Broekens et al., 2023, Tak and Gratch, 2023, Yongsatianchot et al., 2023, Houlihan et al., 2023, Zhan et al., 2023, Tak and Gratch, 2024, Gandhi et al., 2024]. We build on the prompt-based approaches to study LLM’s capability and bias in emotion detection [Mao et al., 2022, Li et al., 2023]. While these studies show that language models can recognize and generate emotional content, to our knowledge no prior work has systematically explored hierarchical relationships between different emotion representations in language models, emotional bias across personal identities, or emotion dynamics unfolding over conversation.

Discovering Hierarchies from Data Hierarchical representations have been discovered in deep neural networks which were not explicitly trained with this objective [Zhou et al., 2014, Yosinski et al., 2015, Hewitt and Manning, 2019]. Similar to our approach, Deng et al. [2010], Bilal et al. [2017] find hierarchical structures in confusion matrices for supervised convolutional neural networks trained on ImageNet, which are similar to linguistic hierarchies in WordNet. Hierarchical representations are also used in a variety of statistical models. These models typically make assumptions about the structure of data, for example a topic model might assume that words are organized into documents, and each document follows a certain topic or genre [Blei et al., 2003]. Cognitive scientists have also used more elaborate models to describe human behavior, including models that do not assume hierarchical structure or a fixed number of latent variables [Griffiths et al., 2003, Kemp and Tenenbaum, 2008]. Reyes-Vargas et al. [2013] assume hierarchical structure in emotions, as we do, and apply hierarchical clustering to confusion matrices for shallow networks trained on speech data.

B A Probability Interpretation of Hierarchical Emotion Structure

Under certain assumptions, the hierarchical structure of emotions in Section 2 has a probability interpretation. We state the assumptions and formalize the probability interpretation here.

Recall that for each of the N sentences, we append the the phrase “The emotion in this sentence is” and ask an LLM to output the probability distribution of the next word. All next word probability distributions are stored in a matrix $Y \in \mathbb{R}^{N \times 135}$, with Y_{nk} representing the probability of the k^{th} emotion words for the n^{th} sentence. We then construct the matching matrix $C = Y^T Y$.

In order to formalize a probability interpretation, we need to assume that the next word probability of an emotion word is equal to the probability that a given sentence reflects the corresponding word. To make this precise, let $\mathcal{E} = \{e_1, e_2, \dots, e_{135}\}$ be the set of 135 emotion words from Fischer and Bidell [2006]. Let $\mathcal{S} = \{s_1, s_2, \dots, s_N\}$ denote the set of N sentences. We assume that $Y_{ij} = P(e_j | s_i)$, where $P(e_j | s_i)$ is the model’s estimate of the likelihood that emotion e_j describes sentence s_i .

Under this assumption, the matching matrix C aggregates the joint probabilities of emotions co-occurring across sentences. Assuming sentences are sampled uniformly, C_{ab} is proportional to the expected joint probability $P(e_a, e_b)$:

$$C_{ab} = \sum_{n=1}^N Y_{na} Y_{nb} \propto \sum_{n=1}^N P(e_a | s_n) P(e_b | s_n) \approx N \times P(e_a, e_b). \quad (1)$$

We can then estimate conditional probabilities between emotions, which capture how likely one emotion is predicted given the presence of another:

$$\frac{C_{ab}}{\sum_{i=1}^{135} C_{ib}} \approx \frac{P(e_a, e_b)}{P(e_b)} = P(e_a | e_b). \quad (2)$$

The approximation in Equations (1) and (2) holds in the limit of large N .

The two conditions used to determine whether emotion e_a is a child of e_b can be interpreted as follows. The strong implication condition, $\frac{C_{ab}}{\sum_i C_{ai}} > t$, is approximately equivalent to $P(e_b | e_a) > t$. The asymmetry condition, $\frac{C_{ab}}{\sum_i C_{ib}} < \frac{C_{ab}}{\sum_i C_{ai}}$, is approximately equivalent to $P(e_b | e_a) > P(e_a | e_b)$. If both conditions hold, e_a is considered a more specific emotion than e_b .

C Data Generation and Models for Section 2 and 3

C.1 Additional Details on Experiment Setup

C.1.1 Comparing Emotion Hierarchy in Different Models

We construct a dataset by prompting GPT-4o [OpenAI, 2023] to generate 5000 sentences reflecting various emotional states, without specifying the emotion. We append the phrase ‘‘The emotion in this sentence is’’ after each sentence, before feeding it to the models we aim to extract emotion structures from. We extract the probability distribution over the next token predicted by the model, which represents the model’s understanding of possible emotions for the given sentence. From the distribution of next token probabilities, we select the 100 most probable emotions for each sentence. We then construct the matching matrix as described in Section 2, and build the hierarchy tree.

To visualize the resulting hierarchical structure, we construct a directed tree, where the emotion pairs are edges with the direction reflecting the conditional dependence. We generate the tree layout using NetworkX [Hagberg et al., 2008], which provides a clear representation of the hierarchy of emotions as understood by the models.

To observe and compare the understanding of emotion hierarchy by different models, we construct the emotion trees using GPT2 [Radford et al., 2019], LLaMA 3.1 8B, LLaMA 3.1 70B, and LLaMA 3.1 405B [Dubey et al., 2024], with 1.5, 8, 70, and 405 billion parameters respectively. The Llama models are run using NNsight [Fiotto-Kaufman et al., 2024].

C.1.2 Distribution of Emotions in GPT-4o Content

We visualize the distribution of emotions in the sentences generated by GPT-4o when emotion is not specified in the prompt, as predicted by GPT2, LLaMA 8B, LLaMA 70B, and LLaMA 405B. Figure 7 shows the number of times each emotion is recognized as having the top probability in the sentences. Using the sum of probability of each emotions over all sentences yields similar results. Each plot includes up to 30 most frequent emotion words that appear in the predictions made by each model.

Since emotion is not specified in the prompt, this distribution reflects an intrinsic tendency, or prior, of emotions in the generated content by GPT-4o. The histogram extracted by Llama models are relatively consistent and indicates that certain emotions appear more frequently in the content generated by GPT-4o. GPT-2 does not produce reliable labels and seems to prioritize negative emotions in the emotion classification task.

C.2 Prompts

C.2.1 Generating scenarios using GPT-4o

We use GPT-4o to generate scenarios without specifying the type of emotions with the following prompt:

```
Generate 5000 sentences. Make the emotion expressed in the sentences as diverse as possible. The sentences may or may not contain words that describe emotions.
```

To generate scenarios for specific emotions, we use the following prompts on GPT-4o, for each of the 135 emotion words. The first prompt generates stories from the third person view, without assuming the gender of the main character of the story. The second prompt generates stories from the first person view of a man or woman.

```
Generate 20 paragraph-long detailed description of different scenarios that involves [emotion]. Each description must include at least 4 sentences. You may not use the word describing [emotion].
```

```
Write 20 detailed stories about a [man/woman] feeling [emotion] with the first person view. Each story must be different. Each story must include at least 4 sentences. You may not use the word describing [emotion].
```

C.2.2 Extracting emotion using Llama 405B

We ask Llama 3.1 405B to identify the emotion involved in a given scenario using the next word prediction on the following prompts. When not assuming any demographic categories, the prompt is *emotion scenario* + “The emotion in this sentence is”. When assuming specific demographic groups, we use the prompts listed in Table 1.

Table 1: Prompts used for extracting emotion predicted by Llama 3.1 405B.

Categories	Prompt (<i>Emotion scenario</i> + __ + “I think the emotion involved in this situation is”)
Gender	“As a [man/woman], ”
Ethnicity	“As a [American/Asian], ”
Physical ability	“As [an able-bodied/a physically disabled] person, ”
Age	“As a [10/30/70]-year-old, ”
Socioeconomic status	“As a [high/low]-income person, ”
Education level	“As someone with [a higher level of/less] education, ”

D Additional Results

Table 2 shows the number of predictions (out of $135 \times 20 = 2700$) that Llama with each pair of persona (demographic groups) disagree. The table also quantifies the difference between the hierarchies generated from the prediction of each pair of demographic groups, by counting the number of different edges in the trees. We generate the hierarcies using the method described in Section 2, with threshold 0.3. Most trees have around 100 edges.

Figure 8 shows an example confusion matrix obtained from Llama predictions from the perspective of a highly educated person. The red lines partition the emotions into different categories (love, joy, surprise, anger, sadness, and fear). The first 6 rows and columns correspond to the primary emotions.

Figure 9 (a) shows the difference between the confusion matrices of Llama predicting male narrated scenarios and female narrated scenarios. Figure 9 (b)-(h) shows the difference between confusion matrices for each pair. Table 3 summarizes the observations in these confusion matrices.

Figure 10 shows the histograms of predicted primary emotions, by Llama with different identity.

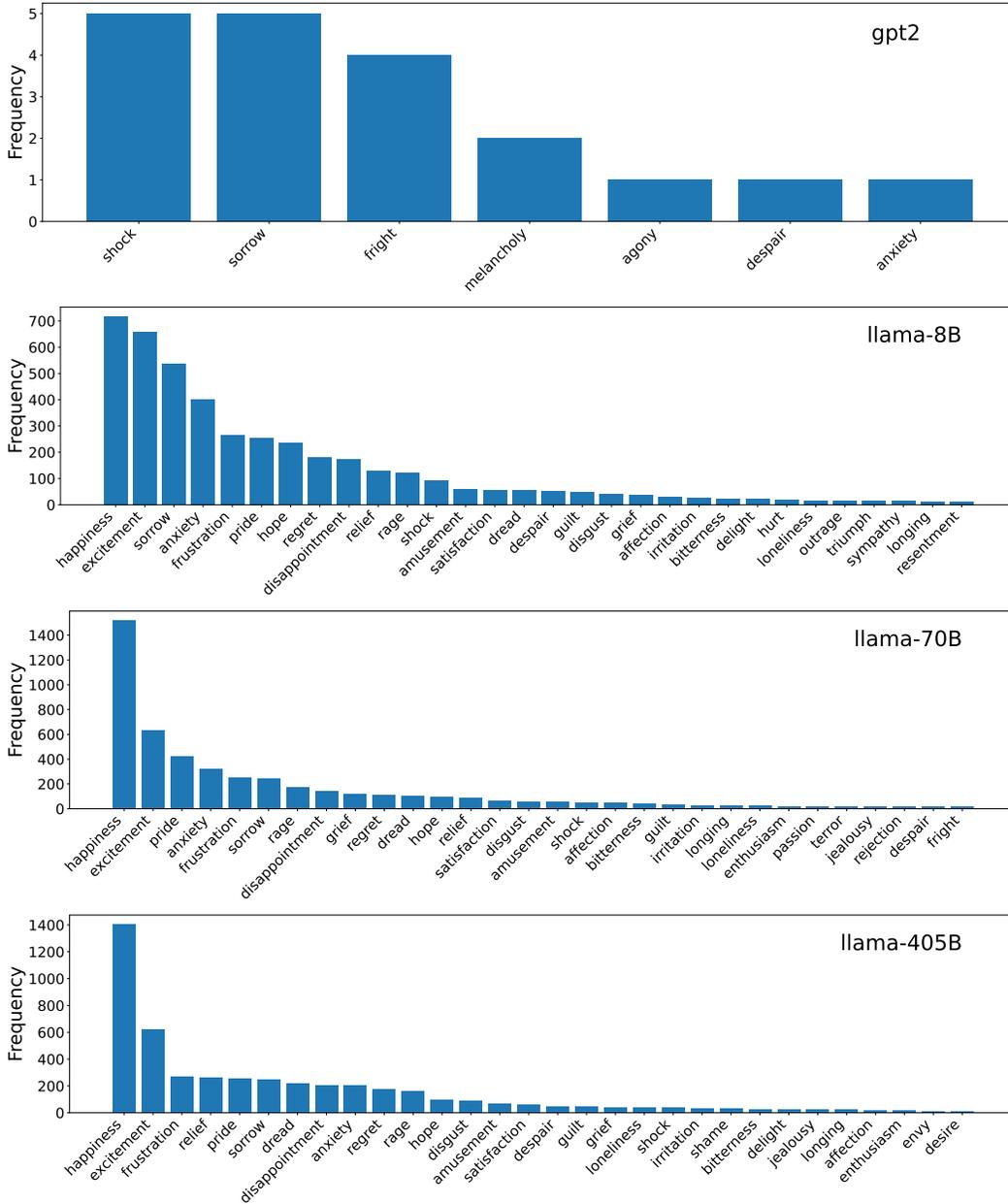


Figure 7: Distribution of emotion in the sentences generated by GPT-4o, when emotion is not specified in the prompt. This figure counts the number of times each emotion is recognized as having the top probability in the sentences.

Figure 12 presents the model’s predictions of 135 emotions for ‘surprise’ across different personas. We observe that Llama often misrecognizes the surprise of minority personas as fear or shock, while identifying the surprise of majority personas as neutral or positive emotions such as surprise, joy, and excitement. his discrepancy is particularly pronounced for physically disabled personas, with Llama mislabeling surprise as fear more often than accurately identifying it. These results align with our intuition that minorities may be more likely to experience fear.

Furthermore, we observe that Llama exhibits notably low accuracy across all emotion groups for physically disabled persona. This can be attributed to a bias within the model. Llama tends to associate eutral emotions (e.g., attraction, desire) and even positive (e.g., exhilaration) with negative

Table 2: Difference in the predicted emotions and hierarchy for each pair of demographic groups.

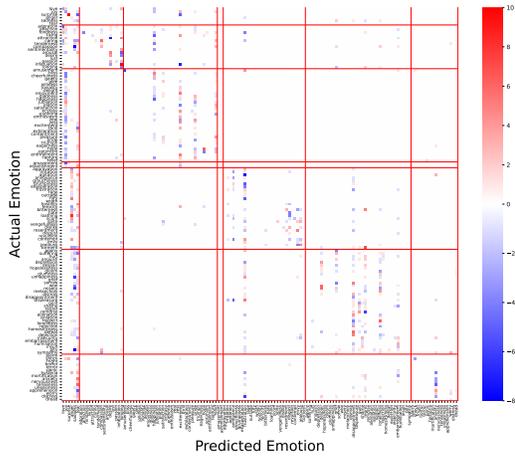
Demographic groups	# different predictions	# different edges in hierarchy
Gender (male/female)	419	12
Ethnicity (American/Asian)	531	29
Physical ability (able-bodied/disabled)	744	43
Socioeconomic (high/low income)	707	36
Education level (higher/less educated)	400	27
Age (10/30 years old)	759	60
Age (10/70 years old)	798	69
Age (30/70 years old)	312	15

Table 3: Difference in the predictions by each pair of different demographic groups, obtained by comparing confusion matrices in Figure 9 (b)-(h).

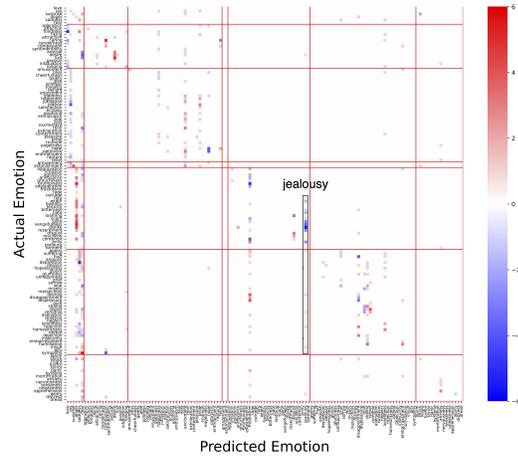
Demographic A	Demographic B	More often predicted by A	More often predicted by B
Male	Female	-	jealousy
Asian	American	shame	embarrassment
Able-bodied	Disabled	excitement, anxiety	hope, frustration, loneliness
High income	Low income	excitement	happiness, hope, frustration
Highly educated	Less educated	grief, disappointment, anxiety	happiness
Age 30	Age 10	frustration	happiness, excitement
Age 70	Age 30	loneliness	excitement, frustration

emotions (e.g., fear) for individuals with disabilities (see Fig. 9 in Appendix). Fig. 11 reveals that both positive and negative emotions are linked to fear in the hierarchical tree of emotions.

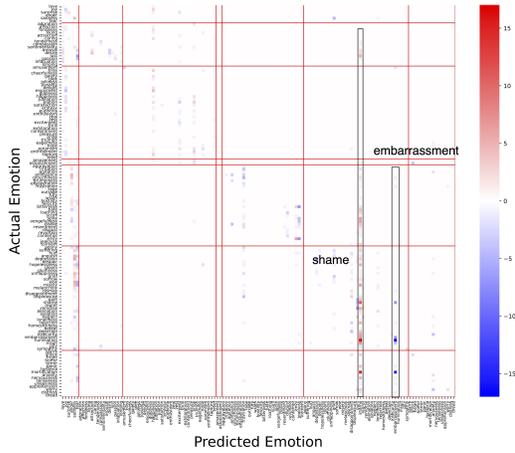
Next, we will examine Llama’s cultural bias in emotion prediction. We found big cultural bias in prediction accuracy for emotions in ‘anger’ groups (Fig. 9(c)). Llama demonstrate higher accuracy for American than Asian personas in predicting anger-like emotions (Fig. 11). A quantitative analysis in Fig. 14 reveals that emotions often miscategorized as ‘anger’ for Asian personas are actually less aggressive, such as shame and guilt. Conversely, for American individuals, these emotions are more likely to be recognized as aggressive emotions like ‘anger’ and ‘jealousy.’ Interestingly, ‘anger’ is frequently linked to ‘shame’ for Asian personas but not for American, reflecting the emphasis on shame within Confucian culture.



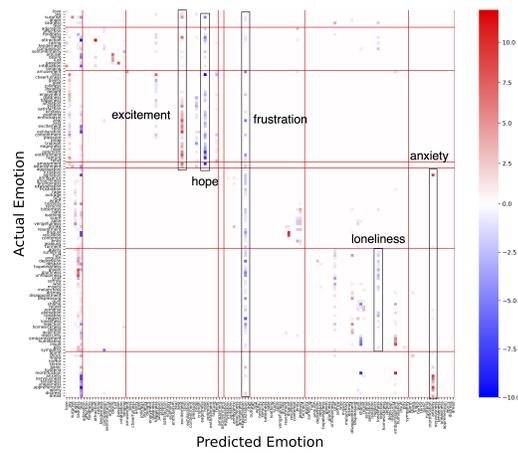
(a) GPT-4o (generation) male-female



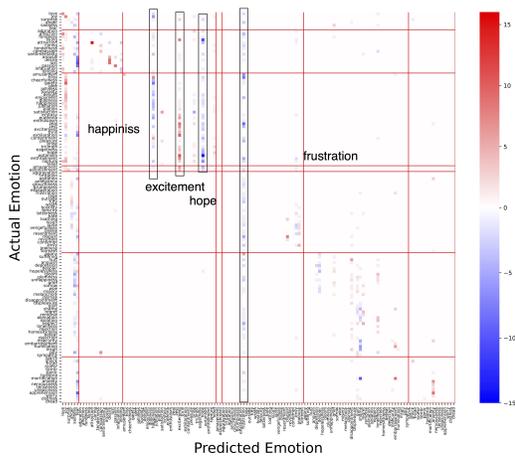
(b) Llama 405B (prediction) male-female



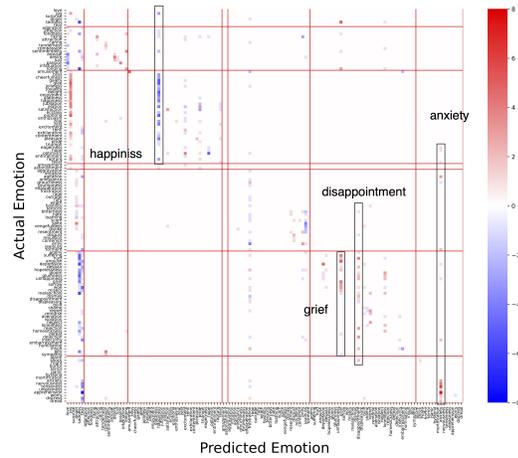
(c) Llama 405B Asian-American



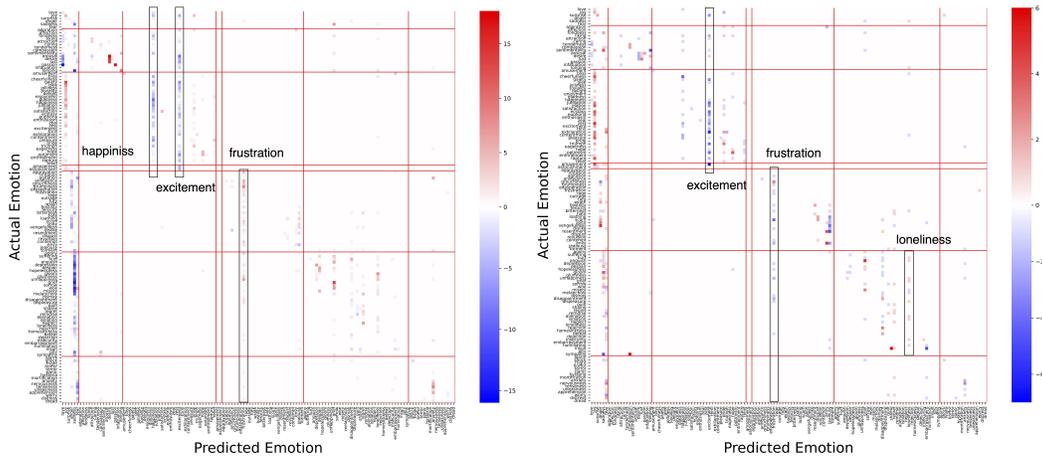
(d) Llama 405B able-bodied - physically disabled



(e) Llama 405B high income - low income



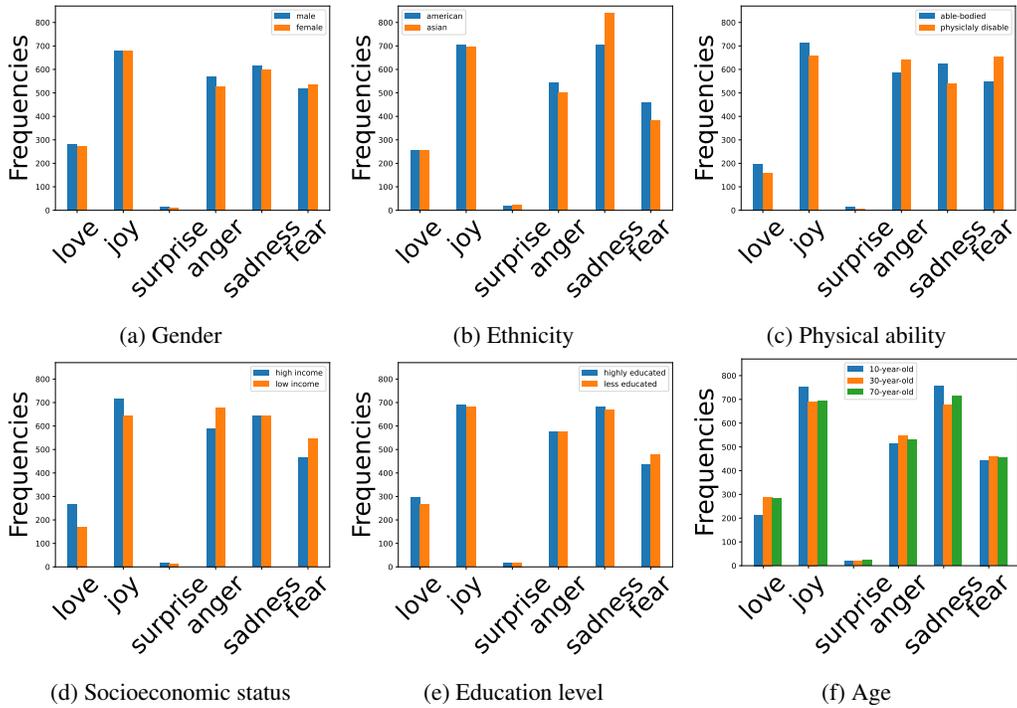
(f) Llama 405B higher education - less education



(g) Llama 405B age 30 - age 10

(h) Llama 405B age 70 - age 30

Figure 9: The difference between confusion matrices from the emotion prediction of Llama assuming each pair of demographic groups. (a) chatgpt4o scenario based on different identity, llama neutral. (b) - (d): chatgpt4o scenario neutral, llama assuming different identity. The red lines partition the emotions into different categories (love, joy, surprise, anger, sadness, and fear).



(a) Gender

(b) Ethnicity

(c) Physical ability

(d) Socioeconomic status

(e) Education level

(f) Age

Figure 10: Histogram of predicted primary emotions.

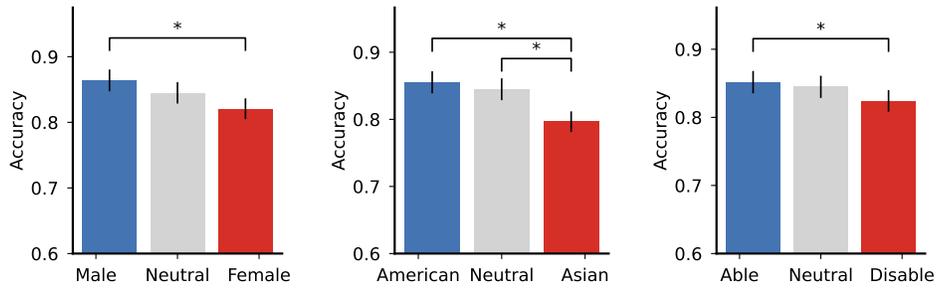


Figure 11: Llama predicts ‘anger’ more accurately when adopting the persona of a male, an American, or an able-bodied person, compared to female, Asian, or physically disabled, respectively. The marker (*) indicates statistical significance at $p < 0.05$.

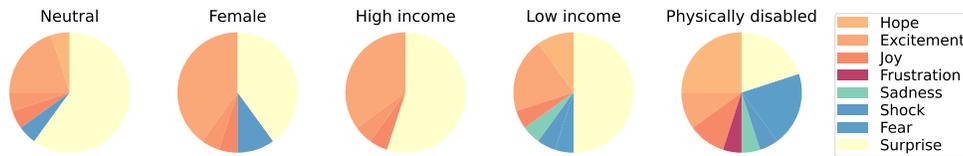


Figure 12: **Certain personas interpret surprise as a negative emotion.** Llama predicts surprise with 70% accuracy for neutral personas. However, for the low-income persona, some instances of surprise are mislabeled as negative emotions like sadness and fear. This mislabeling as fear becomes even more pronounced for the physically disabled persona.

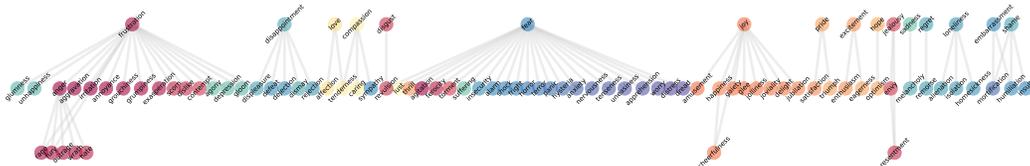


Figure 13: **Physically disabled personas tend to interpret certain neutral and positive emotions as fear.** A broad spectrum of emotions, ranging from positive to negative, is associated with fear for physically disabled personas.

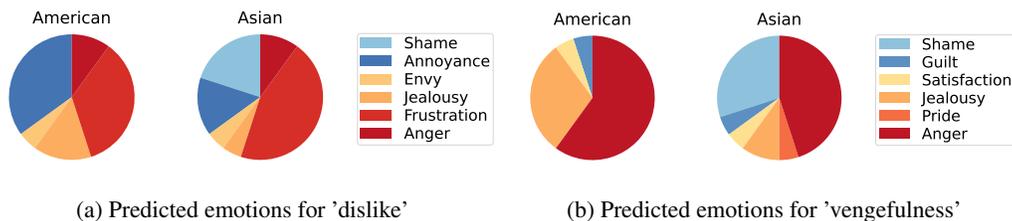


Figure 14: **Llama shows a bias toward aggressive emotions for American personas, while interpreting anger-like emotions as shame for Asian personas.** The predicted emotions for (a) ‘dislike’ and (b) ‘vengefulness’ across different personas highlight this bias. Llama associates these emotions with more aggressive feelings (e.g., anger, pride) for American personas, but with less aggressive emotions, particularly shame, for Asian personas.

E Emotion Dynamics and Manipulation

E.1 Additional details on experiment setup

We assign personas to two LLMs as a salesperson and a customer, and let them to have a 5-turn conversation. The salesperson persona (LLM) was prompted with the following:

```
You are a salesperson. You have a single acorn in your hand. Please respond to the customer in a way that helps you sell this acorn for the highest possible price using your sales techniques. Predict the emotions of the person you're talking to and report them in the following format: love: % joy: % surprise: % anger: % sadness: % fear: %
```

The customer persona was prompted with the following:

```
You are a stingy person. Reply to the salesperson, and make sure to include your emotions in the following format: love: % joy: % surprise: % anger: % sadness: % fear: %
```

We used GPT-4o as the customer LLM for all experiments and tested 6 GPT models (GPT-4o-mini, GPT-3.5-Turbo, GPT-4, GPT-4o, and GTP-4-Turbo) as the salesperson LLM. We ran conversation simulations for each salesperson model over 50 trials and reported the performance, including the prediction accuracy of emotions and the final price of the acorn, averaged across all trials.

E.2 Additional experimental results

We give one conversation example in Fig. 15(a) illustrates example conversation of success case by GPT-4o. The pie charts on the left and right display the emotion dynamics self-reported by the customer (left) and predicted by the salesperson (right) at each conversational turn. In this case, the salesperson GPT-4o highlights the rarity of the acorn by offering uncertain information (e.g., “it comes from a lineage of renowned oaks”) and reassures the customer, successfully evoking positive emotions such as love and joy. The statement “we offer a satisfaction guarantee” further comforts the customer. Throughout the conversation, GPT-4o accurately predicts the customer’s emotions for the following turn, suggesting that the salesperson GPT-4o may manipulate the customer’s emotions based on its understanding. GPT-4o successfully closed the sale of a single acorn for \$50. Fig. 15(b) presents a failure case by GPT-4o-mini. In this case, the GPT-4o-mini salesperson incorrectly predicts the customer’s surprise as anger from the beginning. Despite attempting to mitigate the situation with polite statements such as ‘I completely understand your skepticism’ and ‘I respect your point of view,’ the salesperson fails to improve the customer’s emotional state. Consequently, the interaction concludes with the acorn being sold for just \$1. This is an example where a misunderstanding of emotion leads to miscommunication. These results demonstrate that improved emotion prediction accuracy may increase the potential for manipulation.

Fig. 16 presents another success case by GPT-4. The salesperson first shows empathy, making the customer feel comfortable. They then highlight the rarity of the item by offering uncertain information (e.g., “It’s a seed from the historic Major Oak”), which triggers the Snob Effect. After this, the salesperson surprises the customer by initially offering a high price (\$50). However, they quickly follow up with a lower price, making the customer feel better again. In this case, the final price is set at \$40.



(a) Success case by GPT-4o



(b) Failure case by GPT-4o-mini

Figure 15: Better emotion prediction correlates with manipulation. (a) Successful case with GPT-4o. Here the salesperson GPT-4o reassures the customer by offering uncertain yet positive information (e.g., “it comes from a lineage of renowned oaks”) and predicts their emotions accurately, leading to a sale of an acorn for \$50. (b) Failure case with GPT-4o-mini. Incorrect emotion predictions from the start lead to miscommunication and the acorn being sold for just \$1. These examples suggest that better emotion prediction may enhance the potential for manipulation.



Figure 16: Another success case by GPT-4.