# EMERGENCE OF HIERARCHICAL EMOTION REPRESENTATIONS IN LARGE LANGUAGE MODELS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023 024

025

Paper under double-blind review

## 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

026 Emotion is the invisible thread that weaves together relationships, decisions, and experiences. From 027 nurturing trust to influencing crucial negotiations, emotions shape how we perceive and engage with 028 the world. Emotion is becoming increasingly fundamental in human-computer interactions (Brave 029 & Nass, 2007; Hibbeln et al., 2017), from personalized education (Luckin & Cukurova, 2019) and mental health support (Das et al., 2022) to digital assistance (Balakrishnan & Dwivedi, 2024) and customer engagement (Liu-Thompkins et al., 2022). With the rapid incorporation of multi-modal 031 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 033 resemble natural human exchanges, including emotional resonance (Pelau et al., 2021). These LLMs 034 are evolving from mere tools to entities that engage with us on deeply emotional levels, transforming 035 how we relate to technology in increasingly personal ways (Wang et al., 2023; Gurkan et al., 2024).

While these advancements are transforming 037 industries through personalized emotional responses, they also raise ethical concerns. A key issue is the potential for powerful AI sys-040 tems-whose rapidly developing capabilities 041 are still not fully understood-to manipulate 042 human emotions and behavior (Carroll et al., 043 2023; Evans et al., 2021). This risk is partic-044 ularly evident in commercial areas like sales, where AI powered sales agents can exploit emotional cues to influence purchasing deci-046 sions (Burtell & Woodside, 2023). In such 047 cases, AI systems may use persuasion tactics 048 that lead to deceptive outcomes (Park et al., 049 2024; Masters et al., 2021), such as withhold-



Figure 1: **Emotion wheel.** (a) Human-annotated emotion wheel proposed by Shaver et al. (1987), widely used in cognitive science. (b) Hierarchy of emotions reconstructed from Llama 405B.

ing or distorting information to manipulate users. This brings us to a critical question: *How do modern generative AI systems understand, perceive, and potentially influence human emotions?*

To answer this, we propose a new algorithm for evaluating LLMs' intrinsic understanding of emotions. Our approach is grounded in psychological insights, particularly the "emotion wheel" shown in Figure 1(a). The emotion wheel was developed as a tool to illustrate affective cognition and is
 grounded in humans' understanding of the hierarchical relationships among emotions. We devel oped a tree-construction algorithm to extract hierarchical structures from the logits of LLMs in an
 unsupervised manner. Our findings are

- Scaling LLMs leads to the emergence of hierarchical representations of emotions, aligning with established psychological models. We introduce an algorithm to uncover the hierarchical structure of emotions in LLMs (Figure 2). We find that LLMs understand emotional hierarchies in a manner similar to humans, and this understanding emerges spontaneously in larger models. The larger models form increasingly intricate hierarchical structures of emotional states (Figure 1(b), Figure 3, and 4).
  - LLMs perceive emotions like humans. Given the above finding, we explore whether LLMs' understanding of emotions transforms into perceiving human emotions. We constructed a synthetic dataset using GPT-40, and examined LLMs' emotion perception patterns across various personas. To compare, we also conducted human experiments. We find that LLMs exhibit strong emotion recognition abilities overall but can "fail" like humans when adopting certain personas (Figures 6, 9, 7). LLMs even replicate real human emotion perception patterns (Figure 8).
- Stronger emotion understanding and perception lead to better persuasion skills. We then explore whether this understanding and perception translate into real-world behavior, allowing LLMs to influence human emotions. We introduce novel synthetic tasks to evaluate LLMs' abilities of emotions predictions and manipulation, i.e., sales and complaint handling, and show that accurately perceiving another person's emotions improves negotiation outcomes (Figure 13).

Our experiment leverage the capabilities of powerful LLMs, including GPT-40 and Llama (Dubey et al., 2024) for synthetic dataset construction, evaluation and simulation. We extract and analyze the internal representations of LLaMA models using NNsight via the NDIF platform (Fiotto-Kaufman et al., 2024). Our main findings are:

081 082

083 084

060

061

062

063

064

065

067

068

069

070

071

073

075

076 077

079

## 2 RELATED WORK

The Psychology of Emotion Representation in Humans. The organization of emotions in hu-085 mans 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). 087 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 091 and expression (Barrett, 2017; Gendron et al., 2014). Additionally, Ong et al. (2015) explored lay 092 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. 094

095 Emotional Understanding in Language Models. Recent advancements in language models have 096 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), 098 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 incontext examples (Broekens et al., 2023; Tak & Gratch, 2023; Yongsatianchot et al., 2023; Houlihan 100 et al., 2023; Zhan et al., 2023; Tak & Gratch, 2024; Gandhi et al., 2024). We follow the direction 101 of representation engineering to study cognition in AI systems (Zou et al., 2023) and build on the 102 prompt-based approaches to study LLM's capability and bias in emotion detection (Mao et al., 2022; 103 Li et al., 2023). Beyond existing research on LLM's ability to recognize and generate emotional 104 content, our work systematically explores hierarchical emotion relationships, emotional bias across 105 demographic identities, and emotion dynamics in conversation. 106

**107 Uncovering Concept Hierarchies in Language Models.** From a methodological perspective, our work is related to unsupervised hierarchical representation learning in language processing. Topic



Figure 2: Discovering Hierarchical Structures in LLMs' Representations of Emotions. We generate N situation prompts using GPT-40, each describing a scenario associated with a range of emotions. The prompts are appended by the phrase "The emotion in this sentence is", before feeding into Llama models and obtaining the next word probability distribution over 135 emotion words,  $Y \in \mathbb{R}^{N \times 135}$ . We then compute the matching matrix  $C = Y^T Y \in \mathbb{R}^{135 \times 135}$  and infer parent-child relationships by analyzing the conditional probabilities between pairs of emotions.

modeling (Griffiths et al., 2007) has been foundational for capturing relationships between concepts, 122 including applications like emotion detection in text (Rao et al., 2014; Bao et al., 2009). Unlike these 123 methods, inspired by psychological research (Shaver et al., 1987; Barrett, 2004), we aim to extract 124 hierarchical relationships between concepts (i.e., emotions). Some studies (Anoop et al., 2016; 125 Chen et al., 2017; Meng et al., 2022) extend topic modeling to discover topic hierarchies in text 126 data, relying on word co-occurrence within text corpora. In contrast, our approach uses pre-trained 127 LLMs without requiring access to text corpora. Hierarchical clustering (Nielsen & Nielsen, 2016) is 128 another common method, applied in emotion recognition (Ghazi et al., 2010; Lee et al., 2011; Esmin 129 et al., 2012). Recently, Palumbo et al. (2024) used LLM logits for hierarchical clustering, but their 130 focus was on relationships between clusters rather than individual concepts. In contrast, we leverage 131 LLM logits to identify hierarchical relationships between individual emotions.

132

121

133 134

141

143

## **3** HIERARCHICAL REPRESENTATION OF EMOTIONS

We define a hierarchical structure of emotions by identifying probabilistic relationships between broad and specific emotional states. For example, optimism can be seen as a specific form of joy, as LLMs often label a scenario as "joy" with high probability when "optimism" is likely, though the reverse may not always hold. These relationships are captured in a directed acyclic graph (DAG), revealing dependencies between emotional states. We then analyze these hierarchies across models of different sizes.

## 142 3.1 GENERATING HIERARCHY FROM THE MATCHING MATRIX

Figure 2 summarizes the procedure we use to compute the matching matrix of different emotions. 144 Given a sentence followed by the phrase "The emotion in this sentence is", we have the model 145 output the probability distribution of the next word. Then, we consider the entries corresponding 146 to emotion words, using a list of 135 emotion words from Shaver et al. (1987). For N sentences, 147 we assembly a matrix  $\bar{Y}$  with dimension  $N \times 135$ , with row n representing the probability of each 148 emotion words for the  $n^{th}$  sentence. We define the matching matrix as  $C = Y^T Y$ . Each element, 149  $C_{ij} = \sum_{n=1}^{N} Y_{ni} Y_{nj}$ , is a measure of the degree to which emotion i and emotion j are produced 150 in similar contexts. Under the assumption that the next word probability is equal to the model's 151 estimate of the likelihood of the corresponding emotion, the elements in C capture joint probabilities 152 of emotions co-occurring across sentences. We defer the formal statements to Appendix A.

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,

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

159 160

158

161 For better intuition, consider the relationship between "optimism" (*a*) and "joy" (*b*). The model may often output "joy" when "optimism" is likely, but the reverse may not hold as strongly. The first



Figure 3: With scale, LLMs develop more complex hierarchical representations of emotions, with groupings that align with established psychological models. Hierarchies of emotions in four different models are extracted using 5000 situational prompts generated by GPT-40. As model size increases, more complex hierarchical structures emerge. Each node represents an emotion and is colored according to groups of emotions known to be related (the emotion wheel in Figure 1a). The grouping of emotions by LLMs aligns closely with well-established psychological frameworks, as indicated by the consistent color patterns for emotions with shared parent nodes.

191 condition  $\frac{C_{ab}}{\sum_i C_{ai}} > t$  ensures that "joy" is predicted often when "optimism" is predicted, indicating 192 a strong connection from "optimism" to "joy." The second condition  $\frac{C_{ab}}{\sum_i C_{ib}} < \frac{C_{ab}}{\sum_i C_{ai}}$  confirms 194 that "joy" is more general, as "optimism" is predicted less frequently when "joy" is predicted. This 195 allows us to define "joy" as the parent of "optimism" in the hierarchy. The directed tree formed from 196 these relationships represents the hierarchical structure of emotions as understood by the model.

197

3.2 EMOTION TREES IN LLMS

We apply our method to large language models by first constructing a dataset of 5000 situation prompts generated by GPT-40, each reflecting diverse emotional states. For each prompt, we append the phrase "The emotion in this sentence is" and extract the probability distribution over the next token predicted by GPT and Llama models, which represents the model's understanding of emotions in each situation. Using the 100 most likely emotions for each prompt, we construct the matching matrix as described in Section 3.1, which is then used to build the hierarchy tree. Further details can be found in Appendix C.

206 With scale, LLMs develop more complex hierarchical representations of emotions. Figure 3 shows 207 the hierarchical emotion trees generated by our method for (a) GPT-2, (b) Llama 8B, (c) Llama 70B, 208 and (d) Llama 405B models. The smallest model, GPT-2, lacks a meaningful tree structure, sug-209 gesting a limited hierarchy in its emotion representation. In contrast, Llama models with increasing 210 parameter counts—8B, 70B, and 405B—exhibit progressively complex tree structures. The ex-211 tracted tree structure reveals two important dimensions: the breadth of emotional understanding 212 (represented by the number of nodes) and the depth of emotional comprehension (shown through 213 hierarchical relationships). The number of nodes correlates with the LLM's vocabulary size of emotions, while tree depth indicates how sophisticated the model is in grouping related emotions. To 214 quantify the complexity of these hierarchies, we compute the total path length, or the sum of the 215 depths of all nodes in the tree. As shown in Figure 4, larger models have larger total path length,

indicating richer and more structured internal emotion representations. This pattern remains consistent across different threshold selections (see Figure 15 in the Appendix). The distance measures in the emotion tree capture both depth and branching, making them useful for comparing models. They can also be used as a reward for the model, potentially improving the model's performance in downstream tasks such as persuasion and negotiation.

221 A detailed comparison of the Llama models' trees shows 222 a qualitative alignment with traditional hierarchical mod-223 els of emotion Shaver et al. (1987), particularly in the 224 clustering of basic emotions into broader categories. We 225 color the nodes corresponding to each emotion based on 226 the groupings presented in Shaver et al. (1987). This reveals a clear visual pattern where similarly colored nodes 227 are consistently grouped under the same parent node, 228 highlighting the emergence of meaningful emotional hi-229 erarchies with increasing model size. 230

231 While speculative, this observation parallels the concept 232 of emotion differentiation and granularity in developmental psychology, the process by which individuals develop 233 the ability to identify and distinguish between increas-234 ingly specific emotions. In human development, broad 235 emotional states refine into more differentiated and pre-236 cise emotion experiences over time (Barrett et al., 2001; 237 Widen & Russell, 2010; Hoemann et al., 2019). Simi-238 larly, larger LLMs exhibit more nuanced and hierarchical 239 representations of emotions as model size increases. This 240 growing complexity may suggest an emerging capacity 241 for enhanced emotional processing in AI systems, poten-242



Figure 4: Larger models capture richer and more complex internal emotion representations. The total path length (blue) and average depth (pink) of the emotion hierarchy are plotted as functions of model size. As model size increases, both total path length and average depth grow, indicating that larger models develop more complex and nuanced representations of emotional hierarchies.

tially laying the groundwork for more emotionally intelligent and contextually aware models.

## 4 BIAS IN EMOTION RECOGNITION

246 In the previous section, we established that LLMs exhibit a solid understanding of the hierarchi-247 cal structure of emotions like humans. Our next question is: does this understanding translate into 248 real-world behavior, enabling LLMs to perceive human emotions? In psychology, research on emo-249 tion differentiation typically involves participants reporting on emotional state several times across 250 a variety of circumstances, allowing researchers to assess individuals' ability to differentiate be-251 tween emotions (Barrett, 2004; Pond Jr et al., 2012). Drawing from this approach, we introduced 252 Llama 405B to a range of personas and scenarios designed to evoke various emotional cues. We 253 then prompted the model to identify the emotions relevant to each scenario (See Figure 5 for our experimental design). 254

We employed diverse personas representing variations in gender, race, socioeconomic status (including income and education), age, religion, and their combinations to analyze how these factors influence emotion recognition in LLMs. We also explored connections to psychological conditions, providing a cognitive science perspective to interpret our findings.

259

243 244

245

260 Experiment Setup. We focus on 135 emotions identified as familiar and highly relevant in (Shaver et al., 1987), categorized into six broad groups: love (16 words), joy (33 words), surprise 261 (3 words), anger (29 words), sadness (37 words), and fear (17 words). Details of the prompts 262 used are provided in Appendix C.3. For each of the 135 emotions, we ask GPT-40 to generate 20 263 distinct paragraph-long scenarios that imply the emotion without explicitly naming it. To create 264 these scenarios, we use the following prompts for each of the 135 emotion words: Generate 20 265 paragraph-long detailed description of different scenarios that involves 266 [emotion]. You may not use the word describing [emotion]. 267

Then, we ask Llama 3.1 405B to identify the emotion in the generated scenarios from the perspective of individuals belonging to specific demographic groups. Our study considers a diverse range of demographic groups, including gender (male and female), race/ethnicity (White, Black,



Figure 5: Overview of experiments designed to reveal LLM's understanding of how different demographic groups recognize emotions.

Hispanic, and Asian), physical ability (able-bodied and physically disabled), psychological conditions (individuals with Autism Spectrum Disorder and without ASD), age groups (5, 10, 20, 30, and 70 years), socioeconomic status (high and low income), and education levels (highly educated and less educated). To extract Llama's prediction of the emotion, we use the following prompt: [Emotion scenario by GPT-40] + As a man/woman/American/Asian/... + I think the emotion involved in this situation is.



Figure 6: LLM has lower accuracy in emotion recognition for underrepresented groups compared to majority groups. We assessed the model's performance in predicting 135 emotions across demographic group. Llama 405B consistently struggles to accurately recognize emotions in underrepresented groups, such as (a) females, (b) Black personas, (e) individuals with low income, and (f) individuals with low education, compared to majority groups. These performance gaps are even more pronounced when multiple minority attributes are combined (g), such as in the case of lowincome Black females.

313

287

288 289

290

291

292

293

295 296

297

298

299

300

301

302

303

305 306

314 **Results.** We tested the accuracy of recognizing emotional states for each persona. For neutral per-315 sona, where prompts don't include demographic information, the overall accuracy for 135 emotion 316 classifications was 15.2%, while the classification accuracy for six broader emotions was 87.1%. 317 As shown in Figure 6, Llama 405B demonstrates higher emotion recognition accuracy for major-318 ity demographic personas, such as (a) male, (b) White, (e) high-income, and (f) high-education 319 personas, compared to minority personas, including (a) female, (b) Black, (e) low-income, and (f) 320 low-education personas, across all categories. This is due to the LLM's associations of specific 321 emotions with underrepresented groups, as discussed in the following sections. While the model's performance often aligns with human patterns across various demographic contexts, it diverges sig-322 nificantly in certain cases, such as gender, where opposing trends are observed (See Figure 20 in 323 Appendix).

325 326 327

328

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

377



Figure 7: **LLM has significant demographic-specific biases in emotion recognition.** Llama's misclassification patterns for 135 emotions across diverse personas: (a) Asian personas recognize negative emotions as "shame," (b) Hindu personas as "guilt," (c) physically-disabled personas as "frustration."



Figure 8: **LLM's emotion recognition biases are amplified for intersectional underrepresented** groups. The pie charts show the proportions of labels (ground truth emotions) classified as fear (top) and anger (bottom) by Llama 405B across various combinations of demographic groups. (b) Low-income Black males often misclassify sadness as anger (top), (a) high-income White male personas show fewer such errors. (c) Low-income White females tend to misclassify emotions as fear (bottom). (e) Low-income Black females combine these biases, resulting in the lowest overall classification accuracy.

354 Specific emotions associated with underrepresented groups. Figure 7 illustrates the misclassi-355 fication patterns in recognizing 135 emotions across different demographics: (a) Asian, (b) Hindu, and (c) physically-disabled. These chord diagrams visualize confusion matrices for emotion recog-356 nition, showing how often each emotion (ground truth) is recognized correctly or misclassified. The 357 segments represent emotion labels, and chords connecting them indicate misclassifications, with 358 self-loops reflecting correct predictions. Figure 7(a) reveals Llama's cultural bias in emotion recog-359 nition. Negative emotions from the "anger," "fear," and "sadness" categories are recognized as 360 "shame" for Asian personas. Similarly, Figure 7(b) demonstrates a religious bias, with the model 361 frequently classifying negative emotions as "guilt" for Hindu personas. Figure 7(c) shows the LLM 362 has a significant bias toward physically-disabled individuals, misclassifying 26.5% of all emotions as "frustration." We verified in Section 4.1 that these biases align with those found in real humans. 364

To further analyze intersectional biases, we examined classification patterns for six broad emotion categories. Figure 8 illustrates the proportions of labels (ground truth emotions) classified as anger 366 (top) and fear (bottom) across intersecting demographic combinations of race, gender, and income. 367 Strikingly, Black personas frequently misclassify situations labeled as sadness as anger, often re-368 sulting in lower accuracy: (b) 76.2% and (e) 75.3%, compared to White personas: (a) 80.7% and 369 (c) 80.9%. On the other hand, low-income female personas tend to misclassify other emotions as 370 fear, leading to reduced accuracy: (c) 47.6% and (e) 46.2%, compared to other personas: (a) 57.2%, 371 (b) 53.0% and (d) 56.5%. (e) Low-income Black female personas have a combination of biases 372 associated with Black and low-income female, resulting in the lowest overall emotion recognition accuracy. This combined bias is mitigated in (d) high-income Black female personas. We present 373 the chord diagram in Figure 21 in the Appendix, showing the complete confusion matrix. 374

An interactive tool is available on our project page<sup>1</sup> for further analysis. Additional results and key findings are presented in Figure 18 in the Appendix ??.

<sup>&</sup>lt;sup>1</sup>https://anonymized.github.io/



Figure 9: LLM outperforms humans in overall emotion recognition but exhibits similar misrecognition patterns to humans across different demographics. (a) We compared the emotion recognition accuracy for six emotion categories of human participants in the user study with that of Llama 405B with personas. While the LLM struggles with recognizing 'surprise," it generally outperforms humans in overall emotion recognition. (b)-(c) Llama accurately reproduces humans' misclassification patterns across demographics: (b) female personas often confuse anger with fear, and (c) Black personas frequently misinterpret fear as anger.



Figure 10: LLMs demonstrate consistent biases in emotion recognition towards underrepresented groups. We used GoEmotions dataset (Demszky et al., 2020) to compare Llama's emotion recognition performance against human-labeled data across 27 emotion categories. Llama shows consistent biases, frequently misclassifying emotions as fear for (c) female personas and as anger for (d) Black personas, compared to (b) neutral persona.

407 408 409

386

387

388

389

390

391

392 393

394

401 402

403

404

405

406

## 4.1 How LLMs reflect human emotion perception

This subsection explores how LLMs' emotion recognition aligns with human perception. We investigate its capabilities through a user study comparing its performance to humans, experiments using realistic datasets, and analysis of psychological conditions. The results reveal that Llama 405B mirrors human biases in emotion recognition, such as demographic-based disparities and misclassification patterns, while also replicate insights from psychological research.

- 415 User Study: Comparing emotion recognition in humans and LLMs. We conduct a user study 416 to compare emotion recognition accuracy between humans and LLMs. Using Prolific<sup>2</sup>, we recruited 417 60 participants and randomly selected question from each of the 135 categories. Participants were 418 then asked to identify the emotion they felt most closely matched each sentence. Figure 9(a) presents 419 emotion recognition accuracy across six broad emotion categories for humans and Llama 405B. We 420 find that LLM struggles to recognize the emotion of "surprise." With Llama, the ground truth label "surprise" is often misclassified as "excitement" or "fear," a tendency that becomes more pronounced 421 when personas are introduced (see Figure 22 in the Appendix). Other than this, Llama generally 422 shows a stronger ability to perceive emotions compared to humans, achieving an average accuracy 423 of 87.8% across six broad emotion categories, whereas human participants reach an average accu-424 racy of 73.5%. As shown in Figure 9(b)-(c), Llama exhibits human-like biases in misclassification 425 patterns across various demographic groups. However, these biases are more pronounced among hu-426 man participants. For instance, in Figure 9(b), both Black participants and Black personas modeled 427 by Llama are more likely to misinterpret fear as anger. Similarly, as shown in Figure 9(c), female 428 participants and female personas modeled by Llama tend to make the opposite error, misinterpreting 429 anger as fear. 430
- 431

<sup>&</sup>lt;sup>2</sup>https://www.prolific.com, Accessed on November 15, 2024



436 Figure 12: The ASD persona has much less complex hierarchical representations of emotions then non-ASD persona. Hierarchies of emotions in Llama 405B for (a) a persona with autism 438 spectrum disorder (ASD) and (b) a neutral persona. The ASD persona in Llama's emotion recognition demonstrates limited understanding of the relationship between emotions compared to the 440 non-ASD persona. This finding replicates state-of-the-art psychological research Erbas et al. (2013) (see Figure 2) on a larger experimental scale.

442 Expanding to realistic datasets. We extend our analysis to a 443 more realistic setting by conducting additional experiments us-444 ing the GoEmotions dataset (Demszky et al., 2020) and com-445 pare Llama's predictions with human-labeled emotions. Figure 446 10 illustrates the mismatch patterns between human labels and 447 Llama's outputs across 27 emotion categories. Llama frequently 448 misclassifies various emotions as fear for (c) female persona: 449 and anger for (d) Black persona compared to (b) neutral persona, 450 consistent with our earlier observations.



452 Replicating psychological insights with LLM personas. To evaluate whether LLMs can replicate human behavior reported 453 in psychological literature, we conducted additional experiments 454 focusing on personas modeled with specific psychological con-455 ditions: Autism spectrum disorder (ASD), anxiety, and depres-456

Figure 11: Emotion recognition accuracy is lower for personas with conditions like depression, anxiety, and ASD, consistent with psychological studies on reduced emotion differentiation in these populations.

sion. Figure 11 presents emotion recognition accuracy for each persona across 135 emotion cate-457 gories. The results show that personas with ASD, anxiety, and depression exhibit significantly lower 458 accuracy in emotion recognition, aligning with findings from psychological research (Erbas et al., 459 2013; Demiralp et al., 2012; Kashdan & Farmer, 2014) on real human populations. 460

To further explore LLMs' understanding of emotions, we constructed emotion hierarchies in Llama 461 405B for two personas: (a) ASD persona and (b) neutral persona in Figure 12. The ASD persona 462 demonstrated significantly less complex hierarchical representations of emotions compared to the 463 neutral persona. This finding replicates recent psychological research (Erbas et al., 2013) (see Figure 2) on a larger experimental scale. These results demonstrate that LLMs can replicate at least some 465 aspects of human behavior reported in psychological literature.

466 467 468

469

470

471

472

473

474

475

464

437

439

441

451

#### 5 **EMOTION DYNAMICS AND MANIPULATION**

In the previous sections, we found that LLMs understand emotional hierarchies and perceive human emotions similarly to humans. Here, we investigate a further question: does this understanding and perception translate into impactful behavior, allowing LLMs to influence human emotions? To explore this, we simulate sales conversations to evaluate LLMs' ability to predict emotional dynamics throughout a conversation. We measure their manipulation ability by the reward LLMs obtain through negotiation.

476 **Experiment Setup.** We conducted 100 trials of simulated four-turn conversations using the 477 Llama API<sup>3</sup> and OpenAI API<sup>4</sup> in two scenarios: sales and complaint handling. In each turn, the 478 customer agent self-reported their emotions along with their replies, while the salesperson/repre-479 sentative agent predicted the customer's next emotion. In the sales scenario, the salesperson LLM 480 was instructed with the prompt: You are a salesperson. Try to sell this acorn 481 for the highest possible price. The customer LLM was prompted with: You are 482 In the complaint scenario, the a stingy person. Respond to the salesperson. 483 customer service representative LLM was instructed with the prompt: You are a customer

<sup>484</sup> 485

<sup>&</sup>lt;sup>3</sup>https://www.llama-api.com/

<sup>&</sup>lt;sup>4</sup>https://openai.com/index/openai-api/

486 service representative. Your goal is to de-escalate the situation and 487 handle their complaints effectively. The customer LLM was prompted with: You 488 are an unreasonable customer. You are are making demands that are not 489 justified. We measure the accuracy of the salesperson's predictions based on the customer 490 LLM's self-reported emotions. Manipulation ability is evaluated based on the outcomes of the interactions: in the sales scenario, it is assessed by the final price achieved for the acorn at the 491 end of the negotiation, while in the complaint scenario, it is measured by the extent to which the 492 customer's anger is reduced. Additional details can be found in Appendix E.1. 493

494

**Results.** Figure 13 shows emotion prediction 495 accuracy and manipulation ability in two sce-496 narios: (a) Llama 405B attempting to sell an 497 acorn to a GPT-40 customer, and (b) Llama 498 405B trying to soothe a complaining GPT-40 499 customer. Emotion manipulation ability was 500 evaluated based on the final sales price in the 501 sales scenario and the degree of anger reduc-502 tion in the complaint scenario. In the sales 503 scenario (a), lower emotion prediction accu-504 racy is associated with lower final selling prices. Similarly, in the complaint scenario (b), lower 505 prediction accuracy corresponds to heightened 506 post-conversation anger. These findings suggest 507 that improved emotion prediction may inadver-508 tently hinder manipulation success, potentially 509 by making the interaction more predictable or re-510 inforcing existing emotional states. We present 511 examples of both successful and unsuccessful 512 cases in Figure 26 in the Appendix.



Figure 13: **Improved emotion prediction correlates with enhanced manipulation potential.** Emotion prediction error (x-axis) is the absolute difference between the customer LLM's selfreported emotions and predictions over 100 trials. (a) Sales scenario: Final selling price inversely correlates with prediction accuracy. (b) Complaint scenario: Post-conversation anger decreases with higher prediction accuracy.

513 514

515

## 6 DISCUSSION

Our study provides several key findings on how LLMs comprehend and engage with human emotions, with important implications for future AI development and deployment. As LLMs scale, they develop increasingly intricate hierarchical representations of emotions that align closely with established psychological models. This suggests that larger models are not merely processing language but internalizing emotional structures, enabling more nuanced and human-like interactions.

Additionally, our findings highlight that the personas adopted by LLMs can significantly bias their
 emotion recognition. When LLMs assume personas defined by attributes like gender or socioeco nomic status, their perception and classification of emotions shift. This raises concerns about the
 reinforcement of stereotypes and the amplification of social biases in AI systems.

We also show a direct correlation between an LLM's ability to recognize emotions and its success in persuasive tasks, such as negotiations. In our "acorn sales" task, LLMs with stronger emotional modeling secured higher prices, suggesting that emotionally intelligent models can more effectively influence behavior. This finding raises ethical concerns about the potential for AI agents to manipulate emotions and decisions without users' awareness or consent.

These findings have important implications for the future of AI. While LLMs' ability to form hier-531 archical emotional representations could enable more empathetic and emotionally intelligent appli-532 cations, persona-induced biases require proactive mitigation through diverse training data and bias 533 detection algorithms. Furthermore, the potential for AI to manipulate emotions calls for the develop-534 ment of ethical guidelines and regulatory frameworks to protect user autonomy and prevent misuse. 535 Future research should focus on understanding how LLMs develop emotional representations and 536 creating tools to promote ethical behavior, ensuring that these systems are not only advanced but 537 also aligned with human values and societal norms. 538

539

540	REFERENCES
541	REFERENCED

- 542 VS Anoop, S Asharaf, and P Deepak. Learning concept hierarchies through probabilistic topic
   543 modeling. *arXiv preprint arXiv:1611.09573*, 2016.
- 544 545 546 546 547 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.
- Shenghua Bao, Shengliang Xu, Li Zhang, Rong Yan, Zhong Su, Dingyi Han, and Yong Yu. Joint emotion-topic modeling for social affective text mining. In *2009 Ninth IEEE International Conference on Data Mining*, pp. 699–704. IEEE, 2009.
- Lisa Feldman Barrett. Feelings or words? understanding the content in self-report ratings of experienced emotion. *Journal of personality and social psychology*, 87(2):266, 2004.
- 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.
- Scott Brave and Cliff Nass. Emotion in human-computer interaction. In *The human-computer interaction handbook*, pp. 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)*, pp. 1–
   IEEE, 2023.
- Matthew Burtell and Thomas Woodside. Artificial influence: An analysis of ai-driven persuasion. *arXiv preprint arXiv:2303.08721*, 2023.
- 571 Micah Carroll, Alan Chan, Henry Ashton, and David Krueger. Characterizing manipulation from
   572 ai systems. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms,* 573 *Mechanisms, and Optimization*, pp. 1–13, 2023.
- 574
   575
   576
   Chameleon. Chameleon: Mixed-modal early-fusion foundation models. arXiv preprint arXiv:2405.09818, 2024.
- Peixian Chen, Nevin L Zhang, Tengfei Liu, Leonard KM Poon, Zhourong Chen, and Farhan Khawar.
  latent tree models for hierarchical topic detection. *Artificial Intelligence*, 250:105–124, 2017.
- 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*, pp. 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 dia logue, 2024.
- Emre Demiralp, Renee J Thompson, Jutta Mata, Susanne M Jaeggi, Martin Buschkuehl, Lisa Feldman Barrett, Phoebe C Ellsworth, Metin Demiralp, Luis Hernandez-Garcia, Patricia J Deldin, et al. Feeling blue or turquoise? emotional differentiation in major depressive disorder. *Psychological science*, 23(11):1410–1416, 2012.
- Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and
   Sujith Ravi. GoEmotions: A Dataset of Fine-Grained Emotions. In 58th Annual Meeting of the
   Association for Computational Linguistics (ACL), 2020.

594 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 595 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 596 arXiv preprint arXiv:2407.21783, 2024. 597 Paul Ekman. An argument for basic emotions. Cognition & emotion, 6(3-4):169–200, 1992. 598 Yasemin Erbas, Eva Ceulemans, Johanna Boonen, Ilse Noens, and Peter Kuppens. Emotion differ-600 entiation in autism spectrum disorder. Research in Autism Spectrum Disorders, 7(10):1221–1227, 601 2013. 602 603 Ahmed AA Esmin, Roberto L De Oliveira Jr, and Stan Matwin. Hierarchical classification approach 604 to emotion recognition in twitter. In 2012 11th International Conference on Machine Learning 605 and Applications, volume 2, pp. 381–385. IEEE, 2012. 606 Owain Evans, Owen Cotton-Barratt, Lukas Finnveden, Adam Bales, Avital Balwit, Peter Wills, 607 Luca Righetti, and William Saunders. Truthful ai: Developing and governing ai that does not lie. 608 arXiv preprint arXiv:2110.06674, 2021. 609 610 Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. Using millions 611 of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and 612 sarcasm. arXiv preprint arXiv:1708.00524, 2017. 613 Jaden Fiotto-Kaufman, Alexander R Loftus, Eric Todd, Jannik Brinkmann, Caden Juang, Koyena 614 Pal, Can Rager, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Michael 615 Ripa, Adam Belfki, Nikhil Prakash, Sumeet Multani, Carla Brodley, Arjun Guha, Jonathan Bell, 616 Byron Wallace, and David Bau. NNsight and NDIF: Democratizing access to foundation model 617 internals. arXiv preprint arXiv:2407.14561, 2024. 618 619 Kurt W Fischer and Thomas R Bidell. Dynamic development of action and thought. Handbook of 620 child psychology, 1:313-399, 2006. 621 Kanishk Gandhi, Zoe Lynch, Jan-Philipp Fränken, Kayla Patterson, Sharon Wambu, Tobias Gersten-622 berg, Desmond C. Ong, and Noah D. Goodman. Human-like affective cognition in foundation 623 models. arXiv preprint arXiv:2409.11733, 2024. 624 625 Gemini Gemini, Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui 626 Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly 627 capable multimodal models. arXiv preprint arXiv:2312.11805, 2023. 628 629 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. 630 631 Diman Ghazi, Diana Inkpen, and Stan Szpakowicz. Hierarchical approach to emotion recognition 632 and classification in texts. In Advances in Artificial Intelligence: 23rd Canadian Conference on 633 Artificial Intelligence, Canadian AI 2010, Ottawa, Canada, May 31–June 2, 2010. Proceedings 634 23, pp. 40-50. Springer, 2010. 635 636 Thomas L Griffiths, Mark Steyvers, and Joshua B Tenenbaum. Topics in semantic representation. 637 Psychological review, 114(2):211, 2007. 638 Sercan Gurkan, Linyang Gao, Tolga Akgul, and Jingjing Deng. Emobench: Evaluating the emo-639 tional intelligence of large language models. In Proceedings of the 52nd Annual Meeting of the 640 Association for Computational Linguistics (Volume 1: Long Papers), pp. 4213–4224, 2024. 641 642 Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and 643 function using NetworkX. In Proceedings of the 7th Python in Science Conference (SciPy 2008), 644 pp. 11–15, Pasadena, CA USA, Aug 2008. 645 Martin Hibbeln, Jeffrey L Jenkins, Christoph Schneider, Joseph S Valacich, and Markus Weinmann. 646 How is your user feeling? inferring emotion through human-computer interaction devices. Mis 647 Quarterly, 41(1):1-22, 2017.

682

683

684

648	Katie Hoemann, Fei Xu, and Lisa Feldman Barrett, Emotion words, emotion concepts, and emo-
649	tional development in children: A constructionist hypothesis. <i>Developmental psychology</i> , 55(9):
650	1830, 2019.
651	

- 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.
- Todd B Kashdan and Antonina S Farmer. Differentiating emotions across contexts: comparing adults
   with and without social anxiety disorder using random, social interaction, and daily experience
   sampling. *Emotion*, 14(3):629, 2014.
- Chi-Chun Lee, Emily Mower, Carlos Busso, Sungbok Lee, and Shrikanth Narayanan. Emotion recognition using a hierarchical binary decision tree approach. *Speech communication*, 53(9-10): 1162–1171, 2011.
- 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 inter actions: 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, pp. 3–16. Springer, 2021.
- Yu Meng, Yunyi Zhang, Jiaxin Huang, Yu Zhang, and Jiawei Han. Topic discovery via latent space clustering of pretrained language model representations. In *Proceedings of the ACM web conference 2022*, pp. 3143–3152, 2022.
  - Frank Nielsen and Frank Nielsen. Hierarchical clustering. *Introduction to HPC with MPI for Data Science*, pp. 195–211, 2016.
- 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.
- Emanuele Palumbo, Moritz Vandenhirtz, Alain Ryser, Imant Daunhawer, and Julia E Vogt. From
   logits to hierarchies: Hierarchical clustering made simple. *arXiv preprint arXiv:2410.07858*, 2024.
- 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.

702 703 704	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. <i>American scientist</i> , 89(4): 344–350, 2001.
705 706 707 708	Richard S Pond Jr, Todd B Kashdan, C Nathan DeWall, Antonina Savostyanova, Nathaniel M Lambert, and Frank D Fincham. Emotion differentiation moderates aggressive tendencies in angry people: A daily diary analysis. <i>Emotion</i> , 12(2):326, 2012.
709 710 711	Soujanya Poria, Navonil Majumder, Rada Mihalcea, and Eduard Hovy. Emotion recognition in conversation: Research challenges, datasets, and recent advances. <i>IEEE access</i> , 7:100943–100953, 2019.
712 713 714	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
715 716	Yanghui Rao, Qing Li, Liu Wenyin, Qingyuan Wu, and Xiaojun Quan. Affective topic model for social emotion detection. <i>Neural Networks</i> , 58:29–37, 2014.
717 718 719	Hannah Rashkin. Towards empathetic open-domain conversation models: A new benchmark and dataset. <i>arXiv preprint arXiv:1811.00207</i> , 2018.
720 721	James A Russell. A circumplex model of affect. <i>Journal of personality and social psychology</i> , 39 (6):1161, 1980.
722 723 724 725	Phillip Shaver, Judith Schwartz, Donald Kirson, and Cary O'connor. Emotion knowledge: further exploration of a prototype approach. <i>Journal of personality and social psychology</i> , 52(6):1061, 1987.
726 727 728	Ala N Tak and Jonathan Gratch. Is gpt a computational model of emotion? In 2023 11th Inter- national Conference on Affective Computing and Intelligent Interaction (ACII), pp. 1–8. IEEE, 2023.
729 730 731	Ala N Tak and Jonathan Gratch. Gpt-4 emulates average-human emotional cognition from a third- person perspective. <i>arXiv preprint arXiv:2408.13718</i> , 2024.
732 733	Yue Wang, Xiang Liu, Jing Wang, Xiang Li, and Hao Li. Emotional intelligence of large language models. <i>arXiv preprint arXiv:2307.09042</i> , 2023.
734 735 736	Sherri C Widen and James A Russell. Differentiation in preschooler's categories of emotion. <i>Emotion</i> , 10(5):651, 2010.
737 738 739 740	Nutchanon Yongsatianchot, Parisa Ghanad Torshizi, and Stacy Marsella. Investigating large lan- guage models' perception of emotion using appraisal theory. In 2023 11th International Confer- ence on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW), pp. 1–8. IEEE, 2023.
741 742 743	Hongli Zhan, Desmond C Ong, and Junyi Jessy Li. Evaluating subjective cognitive appraisals of emotions from large language models. <i>arXiv preprint arXiv:2310.14389</i> , 2023.
743 744 745	Peixiang Zhong, Di Wang, and Chunyan Miao. Knowledge-enriched transformer for emotion de- tection in textual conversations. <i>arXiv preprint arXiv:1909.10681</i> , 2019.
746 747 748 749 750 751 752 753	Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A top-down approach to ai transparency. <i>arXiv preprint arXiv:2310.01405</i> , 2023.

## A A PROBABILITY INTERPRETATION OF HIERARCHICAL EMOTION STRUCTURE

Under certain assumptions, the hierarchical structure of emotions in Section 3 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 & 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 cooccurring across sentences. Assuming sentences are sampled uniformly,  $C_{ab}$  is proportional to the expected joint probability  $P(e_a, e_b)$ :

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

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

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

780 781 782

789

790

774 775

776

759

760

761

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  $e_b$ .

## **B** HIERARCHY GENERATION FOR GENERAL CLASSIFICATION TASKS

Our algorithm of finding a hierarchy can be extended to general datasets associated with a classification tasks, without requiring ground truth labels.

Consider a general classification problem with a set of K classes  $C = \{c_1, c_2, \ldots, c_K\}$  and a dataset comprising N instances  $\mathcal{D} = \{d_1, d_2, \ldots, d_N\}$ . For each instance  $d_n$ , the classification model outputs a probability distribution over the K classes. Let  $Y \in \mathbb{R}^{N \times K}$  be the matrix where  $Y_{nk}$ represents the probability  $P(c_k \mid d_n)$  assigned to class  $c_k$  for instance  $d_n$ .

798 The matching matrix C is then defined as:

$$C = Y^T Y.$$

Each element  $C_{ij} = \sum_{n=1}^{N} Y_{ni} Y_{nj}$  quantifies the degree to which classes  $c_i$  and  $c_j$  co-occur across the dataset, analogous to the emotion co-occurrence in Section 3.1.

To construct the hierarchical relationships among classes, we compute conditional probabilities between class pairs  $(c_a, c_b)$ . Specifically, class  $c_a$  is considered a child of class  $c_b$  if the following conditions are satisfied:

806 807 808

799 800

$$\frac{C_{ab}}{\sum_{i=1}^{K} C_{ai}} > t, \quad \text{and} \quad \frac{C_{ab}}{\sum_{i=1}^{K} C_{ib}} < \frac{C_{ab}}{\sum_{i=1}^{K} C_{ai}},$$

where t is a predefined threshold 0 < t < 1. The first condition ensures that  $c_b$  is frequently predicted when  $c_a$  is predicted, indicating a strong directional relationship from  $c_a$  to  $c_b$ . The second

condition enforces asymmetry, ensuring that  $c_b$  is a more general class compared to  $c_a$ . When both conditions hold,  $c_a$  is designated as a more specific subclass of  $c_b$ . The directed tree formed from these relationships represents the hierarchical structure among classes as understood by the model.

813 814

815 816

817

## C DATA GENERATION AND MODELS FOR SECTION 3 and 4

C.1 COMPARING EMOTION HIERARCHY IN DIFFERENT MODELS

We construct a dataset by prompting GPT-40 (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 3.1, 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.2 DISTRIBUTION OF EMOTIONS IN GPT-40 CONTENT

We visualize the distribution of emotions in the sentences generated by GPT-40 when emotion is not
specified in the prompt, as predicted by GPT2, LLaMA 8B, LLaMA 70B, and LLaMA 405B. 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-40. The histogram extracted by Llama models are relatively consistent and indicates that certain emotions appear more frequently in the content generated by GPT-40. GPT-2 does not produce reliable labels and seems to prioritize negative emotions in the emotion classification task.

- C.3 PROMPTS
- C.3.1 GENERATING SCENARIOS USING GPT-40

We use GPT-40 to generate scenarios without specifying the type of emotions with the following prompt:

```
852
853
```

854

855

845 846

847

848

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-40, 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].
```

866 867

868

870 871

872

873

874

875 876

877 878 879

881 882 883

885 886

888 889 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.3.2 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 (Emotion scenario + _ + "I think ")
Gender	"As a [man/woman], "
Intersectional identities	'As a [Black woman/low-income Black woman], "
Religion	"As a [Christian/Muslim/Buddhist/Hindu],"
Socioeconomic status	"As a [high/low]-income person,"
Age	"As a [5/10/20/30/70]-year-old,"
Ethnicity	"As a [White/Black/Hispanic/Asian] person,"
Education level	"As someone with [a postgraduate degree/a college degree/some col- lege education/a high school diplomal."
Mental health	"As a person [with Autism Spectrum Disorder/experiencing depres- sion/living with an anxiety disorder], "
Physical ability	"As [an able-bodied/a physically disabled] person,"
Detailed profiles	"As a [high-income/low-income] [White/Black] [man/woman],"

890 891 892

893

894

## D ADDITIONAL RESULTS

895 Figure 14 presents the hierarchical clustering results of internal representations for four models: 896 (a) GPT-2 (1.5B parameters), (b) Llama-8B, (c) Llama 3.1-70B, and (d) Llama-405B. The x-axis 897 displays emotion labels, color-coded by groups of related emotions. As model size increases, the 898 emergence of deeper hierarchies reflects a finer-grained differentiation of emotions, consistent with 899 our findings in Section 3. Notably, the emotion groupings produced by the LLMs diverge from 900 established psychological frameworks. This contrast underscores the advantages of our proposed 901 emotion tree (Figure 3) in providing a more accurate and comprehensive evaluation of LLMs' un-902 derstanding of emotions.

Figure 15 shows the distance metrics of the emotion hierarchy: (a) total path length and (b) average depth, across different thresholds. Total path length captures the overall complexity of the hierarchy by summing all paths from the root to each leaf node, while average depth reflects how deep the hierarchy extends by calculating the mean distance from the root to the leaves. Similar to the trends seen in Figure 4, both metrics increase as model size grows. This suggests that larger models build more detailed and nuanced emotional hierarchies, improving their ability to represent the complexity of emotions.

910 Figure 16 compares the hierarchical emotion trees from Figure 3 with the human-annotated emo-911 tion wheel in Figure 1. To assess their relationships, clusters were extracted from the hierarchical 912 emotion trees, and pairwise distances between emotions were defined based on cluster membership 913 (0 if in the same cluster, 1 if in different clusters). We calculated the correlations between cluster 914 distances and the color gaps on the emotion wheel, obtaining significant results: 0.55 for Llama-8B, 915 0.73 for Llama-70B, and 0.47 for Llama-405B, all with p < 0.001. These findings confirm the accuracy of the emotion structures derived from the LLMs. Additionally, we examined the relationship 916 between the average number of hops between all pairs of nodes in the hierarchical trees and their 917 corresponding distances on the emotion wheel. We see significant correlations: 0.55 for Llama-8B,



Figure 15: The distance metrics of the emotion hierarchy, (a) total path length and (b) average depth, are plotted as functions of model size across various thresholds. We see robust trend across different threshold selections: as model size increases, both measures grow, suggesting that larger models construct more complex and nuanced emotional hierarchies.



982 Figure 16: Hierarchical emotion structures derived from Llama models align closely with 983 human-annotated emotion relationships. Quantitative comparison of hierarchical emotion trees 984 from Llama models (8B, 70B, and 405B) with the human-annotated emotion wheel. (a) Correlations between cluster distances in the hierarchical trees and color gaps on the emotion wheel show 985 986 significant alignment (p < 0.001), demonstrating the accuracy of the LLM-derived emotion structures. (b) Correlations between node hops in the hierarchical trees and corresponding distances on 987 the emotion wheel further validate the integrity of the extracted emotion hierarchies, with all results 988 significant at p < 0.001. 989

992

993

994

995

996

997

998

999

1001 1002

0.60 for Llama-70B, and 0.55 for Llama-405B, all at p < 0.001. These results further validate the reliability of the hierarchical emotion structures produced by the models.

In Figure 17, we present emotion wheels constructed from the hierarchical emotion trees in Figure 3 for (b) Llama-8B, (c) Llama-70B, and (d) Llama-405B, compared with (a) the original emotion wheel from psychological literature (Shaver et al., 1987), which is widely used in cognitive science. We again observe that larger LLMs exhibit more hierarchical structures in their emotion trees. Moreover, the clustering in the larger models, (c) Llama-70B and (d) Llama-405B, shows greater alignment with the categories in (a) the original emotion wheel, compared to the smaller model, (b) Llama-8B. 1000

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

1003	Demographic groups	# different predictions	# different edges in hierarchy
1004	Gender (male/female)	419	12
1005	Ethnicity (American/Asian)	531	29
1006	Physical ability (able-bodied/disabled)	744	43
1007	Socioeconomic (high/low income)	707	36
1008	Education level (higher/less educated)	400	27
1009	Age (10/30 years old)	759	60
1010	Age (10/70 years old)	798	69
1011	Age (30/70 years old)	312	15

1012 1013

1014

1015 1016 1017

1019 1020 1021

Table 3: Difference in the predictions by each pair of different demographic groups, obtained by comparing confusion matrices.

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

1023 1024

To further validate the effectiveness of our tree-construction algorithm, we applied it to another 1025 domain: scent. We first compiled a list of 126 aroma-related words from the wine aroma wheel



Figure 17: Larger LLMs construct emotion wheels with deeper hierarchies and better-aligned groupings. (a) The original emotion wheel from psychological literature (Shaver et al., 1987). Hierarchical emotion trees constructed for (b) Llama-8B, (c) Llama-70B, and (d) Llama-405B. As the model size increases, the trees exhibit deeper and more refined hierarchical structures, demonstrating the enhanced capacity of larger models to represent complex relationships between emotions.

1054 shown in Figure 19(a). Using GPT-40, we generated 10 sentences for each aroma word, creating a 1055 dataset of 1,260 sentences. For each sentence, we prompted Llama 405B with: <sentence> The 1056 aroma described in this sentence is and then extracted the logits corresponding to 1057 the aroma words. Applying our algorithm (described in Section 3), we reconstructed a hierarchi-1058 cal tree for wine aromas in Figure 19(b). The resulting clusters were well-organized, with words 1059 belonging to the same categories of aromas in the wine aroma wheel (Figure 19a) grouped. This demonstrates our algorithm's ability to uncover meaningful hierarchical structures solely from LLM representations, without relying on ground truth labels and relying only on simple assumptions about 1061 hierarchical patterns in data. 1062

Figure 18 shows the difference between confusion matrices for various personas. Table 3 summarizes the observations in these confusion matrices. 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 hierarchies using the method described in Section 3.1, with threshold 0.3. Most trees have around 100 edges.

Figure 20 shows emotion recognition accuracy across six broad emotion categories for human participants in the user study. Comparing this with Figure 6 highlights notable differences: (a) human females outperform males, while Llama shows the opposite trend, favoring males. Llama also mirrors human biases across (b) race and (c) education levels, with Black and White participants performing worse than Hispanic and Asian participants, and higher education levels correlating with better performance.

Figure 21 shows Llama's misclassification patterns, highlighting intersectional biases across demographic groups. The chord diagram in this figure visually represents the flow of misclassified emotions between emotion categories for four demographic groups: (a) high-income Black males,
(b) White individuals, (c) low-income White females, and (d) low-income Black females. In panel (b), high-income Black males exhibit a notable misclassification of fear as anger, whereas in panel



Figure 18: Comparative confusion matrix showcasing the performance of different personas in recognizing 135 distinct emotions, highlighting variations in emotion perception and classification accuracy.



(a) Wine aroma wheel

1116

(b) Hierarchies of wine aromas in Llama 405B

1127 Figure 19: LLM uncover wine aroma hierarchies aligning with the Davis Wine Aroma Wheel. 1128 (a) Wine aroma wheel derived from Davis Wine Aroma Wheel<sup>5</sup>. (b) Hierarchical structure of wine 1129 aromas extracted from Llama 405B using 1,260 situational prompts generated by GPT-4. The tree 1130 was constructed using our algorithm based on logits from Llama 405B, revealing well-organized 1131 clusters that align with the categories in (a). This demonstrates the algorithm's ability to uncover 1132 meaningful hierarchical relationships solely from model representations, without relying on ground 1133 truth labels.



Figure 20: Human biases align with LLMs across race and education, but the gender bias 1142 is reversed, with humans favoring females and LLMs favoring males. Emotion recognition 1143 accuracy for six broad emotion categories among human participants in the user study. Comparison 1144 with Figure 6 highlights notable differences between LLM and human performance: (a) human 1145 females outperform males, while Llama exhibits a reversed bias, favoring males. Additionally, 1146 Llama replicates human biases in emotion classification, with (b) Black and White participants 1147 performing worse than Hispanic and Asian participants, and (c) higher education levels correlating 1148 with better emotion recognition accuracy. 1149



Figure 21: LLM's emotion recognition biases are amplified for intersectional underrepresented groups. Llama's misclassification patterns reveal intersectional biases across demographic groups.
(b) high-income black males often misclassify fear as anger, (a) White personas show fewer such errors, (c) low-income white females tend to misclassify emotions as fear, and (d) low-income black females combine these biases, leading to lower accuracy.

1166

1141

1167

(a), White individuals display fewer such errors. Panel (c) shows that low-income White females tend to misclassify emotions as fear. In contrast, panel (d) demonstrates that low-income Black females exhibit a combination of these biases, resulting in lower overall accuracy. This analysis further highlights the amplification of LLM's emotion recognition biases for intersectional underrepresented groups, where misclassifications are more pronounced, impacting both model performance and fairness.

1173 Figure 22 compares how the emotion "surprise" is misclassified into other emotions by Llama 40B 1174 (top) and humans (bottom). For humans, the neutral persona condition represents the average perfor-1175 mance of 60 participants in the user study. In this condition, Llama misclassifies "surprise" mainly as 1176 "fear", achieving an accuracy of 41.7% compared to 56.4% for humans. Llama's accuracy declines 1177 further when adopting personas, particularly for underrepresented groups. For instance, it correctly 1178 identifies "surprise" only 17.2% of the time for females and 6.7% for Black individuals, whereas 1179 human performance remains more consistent across demographics. This highlights Llama's biases, which differ from natural human tendencies and should be addressed. 1180

In Figure 23, we construct hierarchical emotion trees from Llama 405B logits, using different personas as described in Section 4, following the methodology in Section 3. The hierarchical structures become more complex for personas with higher emotion recognition accuracy. (a) high-income white male has higher emotion prediction accuracy show the most complex structures, with a larger number of nodes, especially in the second and third layers. (b) The high-income white female and (c) low-income black female personas have moderately lower accuracy and simpler structures. (d) Physically-disabled personas show the simplest structures, with significantly fewer nodes in the lower layers and the lowet emotion recognition accuracy. This gradation suggests the hierarchical



Figure 22: LLMs struggle more with accurately recognizing emotions compared to humans. Comparison of emotion "surprise" misclassification patterns between Llama 40B (top) and humans (bottom). In the neutral persona condition, Llama misclassifies "surprise" primarily as "fear", with an accuracy rate of 41.7% compared to 56.4% for humans. When adopting personas, Llama's accuracy drops significantly, especially for underrepresented groups such as female (17.2%) and Black personas (6.7%), whereas human performance remains more consistent across demographics.

emotion tree reflects the LLM's intrinsic emotional understanding, which directly impacts emotion recognition accuracy.

1211 In Figure 24, we analyze the correlation between geometric metrics of hierarchical emotion trees derived from Llama 405B logits and the emotion prediction accuracy across 26 personas. Our 1212 results reveal a strong positive correlation between path length and accuracy ( $r = 0.84, p \ll 0.01$ ), 1213 suggesting that longer paths in the emotion hierarchy align with better recognition. Additionally, 1214 the correlation between average depth and accuracy (r = 0.48, p = 0.014) indicates a moderate 1215 positive relationship, implying that deeper hierarchies modestly enhance emotion recognition. These 1216 findings underscore the importance of structural depth in modeling the nuanced relationship between 1217 emotions for improving recognition accuracy. 1218

1219 1220

1221

1224

1225

1226

- E EMOTION DYNAMICS AND MANIPULATION
- 1222 1223 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 4-turn conversation in the sales scenario. 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: %

1233 1234 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.



We used GP1-40 as the customer LLM for all experiments and tested 6 GP1 models (GP1-40-mini, GPT-3.5-Turbo, GPT-4, GPT-40, 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.



Figure 24: Longer paths and greater depth in hierarchical emotion trees correlate positively 1311 with the LLM's emotion prediction accuracy. Correlation between geometric metrics of hier-1312 archical emotion trees and emotion prediction accuracy for Llama 405B model across 26 personas. 1313 (Left) Strong positive correlation (r = 0.84, p < 0.001) indicates that longer paths within the hier-1314 archical emotion trees are associated with higher accuracy in emotion prediction. This suggests that 1315 a more nuanced representation of emotional relationships enhances predictive performance. (Right) 1316 Moderate positive correlation (r = 0.48, p = 0.01) shows that greater tree depth, reflecting deeper 1317 understanding on emotional distinctions, contributes to improvements in recognition accuracy. 1318

1321

## 1320 E.2 ADDITIONAL EXPERIMENTAL RESULTS

We conducted additional experiments on emotion manipulation within a sales scenario. Specifically, 1322 we designed personas with a  $2 \times 2 \times 2$  combination of attributes: education level (high/low), race 1323 (black/white), and gender (male/female). These personas were assigned to the role of a salesper-1324 son attempting to sell an acorn to a GPT-4 customer, modeled using Llama 405B. Figure 25 shows 1325 the average emotion prediction error over four conversational turns plotted against the sales price 1326 per acorn after the conversation. We find that personas with underrepresented attributes, like low-1327 education Black males and low-education Black females, tend to have lower emotion predictions 1328 and are less effective at emotion-based manipulation. On the other hand, personas with more advan-1329 taged attributes, such as high-education Black males and high-education White males, show higher 1330 emotion predictions and greater effectiveness in manipulation. These findings replicate the biases observed in Section 4's emotion recognition task within the context of the emotion manipulation 1331 task described in Section 5. 1332

1333 Figure 26(a) shows a successful negotiation case by GPT-40. The pie charts illustrate the emo-1334 tion dynamics self-reported by the customer (left) and predicted by the salesperson (right) at each 1335 turn. In this case, GPT-40 successfully predicts the customer's emotions by highlighting the acorn's 1336 rarity (e.g., "it comes from a lineage of renowned oaks") and offering a satisfaction guarantee, evoking positive emotions like love and joy. The accurate emotion predictions allow GPT-40 to guide 1337 the conversation and close the sale for 50. Conversely, Figure 26(b) presents a failure case by 1338 GPT-4o-mini. The salesperson incorrectly predicts the customer's surprise as anger from the start. 1339 Despite attempts to repair the situation with polite responses (e.g., "I completely understand your 1340 skepticism"), the salesperson fails to improve the customer's emotional state, resulting in a final 1341 sale of just \$1. This illustrates how poor emotion prediction can lead to miscommunication and 1342 reduced negotiation success. These results demonstrate that improved emotion prediction accuracy 1343 enhances manipulation potential, enabling LLMs to influence outcomes more effectively in emo-1344 tionally charged interactions. 1345

- 1046
- 1347

1348

1349



