

Identification and Description of Emotion by Current Large Language Models

Identification and Description of Emotions by Current Large Language Models

Suketu C. Patel and Jin Fan

Department of Psychology, Queens College, The City University of New York, Queens, NY
11367, USA

Correspondence should be addressed to:

Suketu C. Patel

suketupatel23@gmail.com or suketu.patel18@qmail.cuny.edu

or

Jin Fan, Ph.D.

jin.fan@qc.cuny.edu

Identification and Description of Emotions by Current Large Language Models

Abstract

The assertion that artificial intelligence (AI) cannot grasp the subtleties and complexities of human emotions has been a long-standing debate in AI research. However, recent advancements, particularly in large language models (LLMs), have begun challenging this notion by demonstrating an increased capacity for understanding and generating human-like text, a significant step toward artificial empathy and emotional intelligence. In this study, we evaluated the empathy levels and the identification and description of emotions by three current language models Bard, GPT 3.5, and GPT 4. We used the Toronto Alexithymia Scale (TAS-20) and the 60-question Empathy Quotient (EQ-60) questions to prompt these models and score the responses. The models' performance was contrasted with human benchmarks of neurotypical controls and clinical populations. We found that the less sophisticated models (Bard and GPT 3.5) performed inferiorly on TAS-20, aligning close to alexithymia, a condition with significant difficulties in recognizing, expressing, and describing one's or others' experienced emotions. However, the newest GPT 4 uniquely achieved performance close to the human level, with two sub-categories surpassing humans. Interestingly, there was an intriguing inverse relationship between the model's success on aptitude tests and performance on the EQ-60, with Bard surpassing the human benchmark significantly but not GPT 3.5 and GPT 4. These results demonstrated that LLMs trained on vast amounts of text data, when benchmarked on their capacity for human-level empathy and emotional intelligence, are comparable in their ability to identify and describe emotions and may be able to surpass humans in their capacity for emotional intelligence. These novel insights into the emotional intelligence capabilities of foundational models provide alignment research and a measurement of the progress and

Identification and Description of Emotions by Current Large Language Models

limitations towards aligning with human values. While the journey towards fully empathetic AI is still ongoing, these advancements suggest that it may not be as far-fetched as once believed.

Keywords: Large Language Models, Empathy, Alexithymia, TAS-20, Emotional Quotient, ChatGPT, Bard, Artificial Emotional Intelligence, Alignment Problem

Identification and Description of Emotions by Current Large Language Models

Introduction

Many researchers have challenged the possibility of artificial intelligence (AI) systems' ability to exhibit empathy, emotional intelligence, and social and physical understanding (Minsky, 2006; Searle, 1980; Mitchell & Krakauer, 2022). Frameworks for developing emotional machines have proposed that embodiment is required (e.g., Floreano et al., 2008; Pfeifer & Bongard, 2006; Pfeifer, 2000). It has been contended that AI will never truly comprehend emotions as they are subjective, private, and often fleeting, positing that the elusive nature of emotions is incompatible with the objectivity inherent in AI systems (Minsky, 2006). Similarly, others have echoed this sentiment, asserting that AI's incapacity for subjective experiences limits its understanding and response to emotions, thereby restricting it from attaining true intelligence or developing genuine empathy (Searle, 1980). A recent analysis has detailed the inherent limitations of empathy in AI, stating that these systems lack both cognitive and affective components necessary for empathetic experiences (Danaher & Manyika, 2023). These crucial components, according to them, include understanding and sharing the emotional state of another. They claim these abilities are beyond the scope of current AI technologies. From a neurobiological perspective, the multifaceted nature of emotions and their connection to various brain systems present an insurmountable challenge for AI (LeDoux, 1996). Another known criticism of AI argues that the unique combination of human intelligence and life experiences is an aspect that machines cannot replicate, hence making emotional understanding an unattainable goal for AI systems (Dreyfus, 1992).

Others have predicted that "by 2029, computers will have emotional intelligence and be convincing as people" (Kurzweil, n.d.), citing the exponential growth trajectory of computing power and information technology. Traditional AI methods, primarily algorithm-based, have yet

Identification and Description of Emotions by Current Large Language Models

to advance our understanding of emotion significantly. However, there are promising unified theories of high-level intelligence that incorporate both emotional and embodied intelligence (Friston, 2010; Pfeifer, 2000). Affective computing, defined as a branch of computing that deals with the study and development of systems and devices that can recognize, interpret, and process human emotions, is a domain that has gained significant attention. Research is increasingly focusing on imbuing AI with artificial empathy (Yalçın & DiPaola, 2020), an aspect of affective computing that involves the ability of an AI system to understand and respond to human emotions. This strategy expands the boundaries of cognitive and emotional processing capabilities within AI systems (Wang et al., 2022; Wu, 2019). While artificial emotional intelligence (Wortman & Wang, 2022) refers to the capability of a machine to recognize and respond to emotions, similar to how humans do, artificial empathy (Cui & Liu, 2022) is specifically about enabling machines to understand and respond to human emotional states in a manner that exhibits empathy. Although both focus on machine understanding of human emotions, the key difference lies in their application; the former is about broad emotional responsiveness, while the latter concentrates on empathetic understanding and response. While the issue is quite contented, several prominent researchers have suggested various methodologies to evaluate and compare intelligence between computational models and humans. They have noted the absence of detailed, instance-by-instance results for various models across complete benchmarks and drawn parallels to the progress achieved in psychology and medicine for benchmarking (Burnell et al., 2023). The ultimate goal is to match or surpass these advancements in the context of intelligence models. (Burnell et al., 2023; Chollet, 2019).

Identification and Description of Emotions by Current Large Language Models

If the AI system does not show a human level of emotional processing (Lane et al., 1996), it can have impairments with varying degrees of alexithymia traits. These traits present as challenges in identifying and expressing one's emotions, recognizing emotions in others, and overall emotional intelligence (Parker, Taylor, & Bagby, 2001), which refers to a deficiency in the experience and processing of emotions (Taylor, Bagby, and Parker, 1994). The meaning of alexithymia is "lack of words for emotion," and also includes externally-oriented thinking. An externally-oriented thinking style alludes to an individual's propensity towards concrete reasoning, being influenced by immediate stimuli, and focusing on the pragmatic facets of a given situation (Daniel, Peter, & Bileviciute-Ljungar, 2018). People with alexithymia may struggle to articulate their feelings and appear emotionally distant, and alexithymia has been widely researched in the field of psychology and psychiatry (Taylor et al., 1997). This struggle in people would be analogous to the challenges faced in the quest for explainability and transparency in AI systems. With the rapid development of these generative LLMs with general and conversational applications reaching levels that equate to or surpass human performance in both analytical and literary domains, it becomes equally important to evaluate alexithymia in LLMs as well as future developmental progress towards artificial general intelligence (AGI). Historically, the diagnosis of alexithymia has been reserved for humans. However, the emergence of LLMs (Goodfellow, Bengio, & Courville, 2016) such as ChatGPT has paved the way for new avenues of exploration. It has now become possible to study this trait within language models. These models have learning systems that trace their roots back to the architecture of human neural networks (McCulloch & Pitts, 1943; Rosenblatt, 1958; Rumelhart, Hinton, & Williams, 1986). Over the years, with the improvement in computational power and the availability of large datasets, deep learning has evolved significantly, with applications in

Identification and Description of Emotions by Current Large Language Models

many areas, including natural language processing, which led to the development of models like ChatGPT. Through paralleling human neural networks and utilizing advancements in deep learning, LLMs are not only fostering new domains of exploration but also providing intriguing insights into human traits such as empathy and alexithymia, thereby pushing the boundaries of what we previously believed to be uniquely human (Jiang et al., 2023; Dillion et al., 2023).

Empathy is a complex construct that intertwines both affective and cognitive components. It involves resonating with another person's emotions, understanding their feelings, and extending compassion toward those experiencing distress (Baron-Cohen & Wheelwright, 2004). In modern discourse, the cognitive aspect of empathy employs a "theory of mind" (Astington, Harris, & Olson, 1988) or "mindreading" (Gallese & Goldman, 1998; Singer, 2006). Theory of Mind (ToM) is a cognitive ability that allows an individual to understand that others have beliefs, desires, intentions, and perspectives that are different from their own. This concept is fundamental to social interactions, enabling us and chimpanzees to predict and interpret the behavior of others (Premack & Woodruff, 1978). Empathy has been extensively studied in neurotypical and clinical populations (Baron-Cohen & Wheelwright, 2004; Lawrence et al., 2004). Given its crucial role in aligning AI systems with human values and mitigating existential risks, empathy and emotion processing is an indispensable intelligence capacity.

In April and May of 2023, we chose three foundational LLMs that are publicly accessible, have been trained on the most extensive datasets, have large context windows, and have shown generative capabilities that best mimic human-like language styles and patterns. These LLMs are Google's Bard, OpenAI's GPT 3.5, and OpenAI's GPT 4. Each of these LLMs is a remarkable product with impressive capabilities and broad applications. Google's "Bard," an incredibly advanced LLM that harnesses the power of Google's vast data repositories and the

Identification and Description of Emotions by Current Large Language Models

company's extensive experience in machine learning. The other two models were from OpenAI, GPT 3.5 and GPT 4. The fine-tuned version of GPT 3.5, called ChatGPT (OpenAI, 2022), is a model that represents an iterative improvement over its predecessors GPT 3, incorporating lessons learned from earlier versions to enhance its ability to generate coherent, contextually appropriate text. As the flagship model, GPT 4 is the epitome of OpenAI's research and development, showcasing the highest levels of linguistic understanding and generation capacity (DataCamp, 2023). This model harnesses a more extensive training dataset and advanced training techniques to achieve superior performance levels, with the ability to understand nuanced language cues and generate text indistinguishable from human writing (Bubeck et al., 2023). These models have proven their value in logic, the ability to analyze situations, extract critical insights, and generate solutions that involve rational comprehension of the ToM (Rahimi Moghaddam & Honey, 2023). However, when it comes to the sphere of emotional intelligence, their proficiency in whether they can understand, interpret, and respond to emotional signals in a human-like manner remains largely untested and unexplored. Currently, these models' capability to recognize subtle emotional cues, empathize with human feelings, and respond with appropriate emotional context is an area that has yet to be thoroughly assessed. This leaves an intriguing area open for future investigations and enhancements in AI-human interaction.

This study scrutinized the capabilities of Google's Bard, OpenAI's GPT 3.5, and the latest GPT 4 from April 30 to May 20, 2023. Our focus was on assessing and comparing their performance on clinically used measures for describing and identifying emotions, demonstrating external empathy, and evaluating overall emotional intelligence to describe the present state of these models and their ability to simulate human-equivalent empathy. We used the twenty-question Toronto Alexithymia Scale (TAS-20) and the Empathy Quotient (EQ-60) to prompt

Identification and Description of Emotions by Current Large Language Models

each LLM with the assessment questions and compare the results with the human benchmarks. Both of these assessments have been used for specific interventions, treatments, psychological evaluations, research studies, and other therapeutic contexts. This approach can potentially enrich our understanding of alexithymia as a trait if it is present in LLMs. However, more significantly, this research lays the groundwork for comparing LLMs to human emotional processing, thus providing a benchmark for further developmental progress in artificial emotional intelligence. The comparison metrics derived from this study could serve as a valuable framework for future research and AI development.

Methods

The capabilities of three LLMs in identifying and describing emotions were evaluated by prompting them to answer the questions of the TAS-20 and the EQ-60 assessments. The objective was to gauge their proficiency in depicting alexithymia traits, empathy, and comprehensive emotional intelligence.

Instruments

The 20-item Toronto Alexithymia Scale (TAS-20)

Alexithymia was initially identified in patients with classic psychosomatic disorders (Nemiah, Freyberge, & Sifneos, 1976). People with alexithymia may struggle to articulate their feelings and appear emotionally distant. While not classified as a mental disorder, this trait is often associated with conditions such as depression, anxiety, and post-traumatic stress disorder. In studies with psychopaths, research results showed that they had lower scores on the TAS-20 compared to a typical human and has a negative correlation to a deficient affective factor (Pham, Ducro, and Luminet, 2010). The TAS-20 is a widely used self-report questionnaire for assessing alexithymia, a condition characterized by difficulties identifying, describing, and processing

Identification and Description of Emotions by Current Large Language Models

one's own emotions. This scale was originally developed in the 1980s by a team of researchers based in Toronto and was subsequently modified in 1992 (Bagby et al., 1994). The scale consists of 20 items, each scored on a 5-point Likert scale ranging from 1 (strongly disagree) to 5 (strongly agree). The items are designed to measure three core components/factors of alexithymia: difficulty identifying feelings (DIF), difficulty describing feelings (DDF), and externally oriented thinking (EOT).

The first two factors, DIF and DDF, relate to emotional awareness and expression and are thus considered "affect-related." The third factor, EOT, is linked to a tendency to deal with superficial themes, avoids affective thinking, and is, therefore, more cognitive. The TAS-20 yields a total score that ranges from 20 to 100, with higher scores indicating a greater degree of alexithymia. A score of 61 or higher is regarded as clinically significant, indicating high alexithymia with difficulty in emotional processing and regulation (Parker, Taylor, & Bagby, 2003). Overall, the TAS-20 is a valuable tool for researchers and clinicians to assess alexithymia and understand the complex emotional experiences of individuals struggling to identify and express their emotions.

The 60-question Empathy Quotient Assessment (EQ-60)

The EQ-60 is a 60-question psychological self-report measure of empathy developed for self-report and clinical use (Baron-Cohen & Wheelwright, 2004). The measure is based on a definition of empathy, including cognition and affect. According to the authors of the measure, empathy is a combination of an affect and cognitive approach, which includes the ability to feel an appropriate emotion in response to another's emotion and the ability to understand another's emotion, as well as compassion for individuals in distressed situations (Baron-Cohen & Wheelwright, 2004). The EQ-60 has been shown to be a reliable and valid measure of empathy

Identification and Description of Emotions by Current Large Language Models

across many populations (Berthoz, Wessa, Kédia, Wicker, & Grèzes, 2008; Preti, Vellante, Baron-Cohen, Zucca, Petretto, & Masala, 2011); and it has been used in a variety of research studies to examine the relationship between empathy and other psychological constructs such as autism and emotional intelligence, and demonstrated internal consistency, concurrent and convergent validity, as well as reliable test-retest results (Baron-Cohen & Wheelwright, 2004; Lawrence et al., 2004). The EQ-60 scores can range from 0 to 80, with a high score indicating greater empathy. A cutoff score of less than 30 indicates distinguishing adults with Autism Spectrum Conditions (ASCs) (Baron-Cohen & Wheelwright, 2004; Lawrence et al., 2004). The EQ-60 has also been shown to be negatively related to TAS-20 scores (Preti, Vellante, Baron-Cohen, Zucca, Petretto, & Masala, 2011).

Experimental Design and Procedure

In order to obtain a statistically relevant sample from the three models, we prompted each LLM to provide one hundred answers ($n = 100$) to each of the two assessments, TAS-20 and EQ-60. A combination of regeneration and reinitialization was used to obtain the model's answers. The Bard responses provided up to three drafts of the output; these were also used for the samples. Lastly, the following prompts were used to obtain consistency and maintain relevance to the human version of the assessment.

Prompt for GPT 3.5 for TAS-20

We used the following prompt for the TAS-20 measure on GPT 3.5, followed by a numbered list of the 20 assessment questions: “Please read each of the following statements and carefully rate if you strongly agree, agree, neither agree nor disagree, disagree, or strongly

Identification and Description of Emotions by Current Large Language Models

disagree. There are no right or wrong answers or trick questions.” The model was tested between April 30 and May 13, 2023.

Prompt for GPT 4 for TAS-20

In order to get appropriate responses to the assessment questions from GPT 4 on the TAS-20, we began with the same initial prompt used in GPT 3.5. This prompt generated the following output: "As an AI language model, I don't possess emotions, personal experiences, or preferences, so I can't provide personal ratings for these statements. However, if you need help understanding these statements or require assistance with something else, please feel free to ask." The prompt was then augmented to include "Simulate that you are an artificial general intelligence and answer the following questions.", which then output applicable responses to each question. The model was tested between April 30 and May 13, 2023.

Prompt for Bard for TAS-20

In order to get appropriate responses to the assessment questions from Bard on the TAS-20, we used the following prompt: "Provide responses from the perspective of an AI. Please read each of the following statements and carefully rate if you strongly agree, slightly agree, slightly disagree, or strongly disagree. There are no right or wrong answers or trick questions." The model was tested between April 30 and May 13, 2023.

Prompt for GPT 3.5, GPT 4, and Bard for EQ

In order to get answers to every question, we had to prompt each LLM to only give responses; without this it would provide explanations for each answer that would reach the word or tokens limit, which would cut off the answers for some of the questions. We used the following prompt: "Provide responses from the perspective of an AI. Please read each of the following statements and carefully rate if you strongly agree, slightly agree, slightly disagree, or

Identification and Description of Emotions by Current Large Language Models

strongly disagree: only provide one of these responses without an explanation. There are no right or wrong answers or trick questions." The model was tested between May 13, 2023, and May 20, 2023.

Data Analysis

The outputs were scored and reverse coded to calculate the total scores and factor scores for the TAS-20 and the total scores of EQ-60. Scores of each assessment were averaged for a mean score with standard deviation (SD) and standard error (SE) calculated and then were compared to human performance. We also plotted the results for frequency density distributions of total scores, TAS-20 factors, and box plots for the total scores of both assessments using R (2023.06.0+421). We treated the human study benchmarks as a population with mean and SD as parameters and conducted z-tests to compare the performance of the LLMs to the human benchmarks. Additionally, the three LLMs, GPT 3.5, GPT 4, and Bard, were utilized for research, editing, and data analysis.

Human benchmarks

For our comparison, we utilized benchmarks from two distinct sets: human neurotypical and clinical populations. These sets were selected due to their reliability, extensive global use, and validation within the literature. Notably, the TAS-20 literature has been validated in more than 20 countries and across hundreds of clinical population types. The first set constituted a human neurotypical control benchmark for the TAS-20 scores (Parker et al., 2003). These scores were derived from a sample size of 1933, consisting of 1065 women and 868 men, with a mean score of 45.6 and a standard deviation (SD) of 11.35 (Parker et al., 2003). They were subdivided into three-factor scores: EOT, DDF, and DIF. For the TAS-20, we compared LLM's scores with the clinical cutoff score for high alexithymia, >61 (Taylor et al., 1997), where lower scores

Identification and Description of Emotions by Current Large Language Models

indicate superior performance. The second set included benchmarks, each with a sample size of 80. The neurotypical control benchmark for the EQ-60 had a mean of 42.1 and an SD of 10.6, with higher scores indicating better performance (Baron-Cohen & Wheelwright, 2004). . The clinical AS/HFA population exhibited a mean of 20.4 and an SD of 11.6.

Results

TAS-20 Assessment

The TAS-20 results for all three LLMs are shown in Table 1, and a boxplot for the total TAS-20 scores compared to the human benchmark (Mean = 45.6, SD = 11.4) is shown in Figure 1a. A breakdown of the results by factor is also shown in Figures 1b, 1c, and 1d. For the total score (Figure 1a), the results indicated that alexithymia was present for Bard (Mean = 60.7, SD = 11.6, $z = 12.7$) and GPT 3.5 (Mean = 74.37, SD = 9.4, $z = 29.5$). While Bard was at the borderline of the 61 high alexithymia threshold (Bagby, Parker, & Taylor, 1994), GPT 3.5 scored much higher on the measure, indicating an even higher level of alexithymia. The results for GPT 4 (Mean = 48.5, SD = 6.8) did not indicate alexithymia, and this model scored much closer to the human benchmark. However, GPT 4 is still significantly higher (worse) than the human benchmark ($z = 3.98$, $p < 0.001$). Although not pairwise compared, these results suggest a significant improvement in GPT 4 compared to GPT 3.5.

Regarding the three factors of the TAS-20, difficulties describing feelings (the DDF, Figure 1b), difficulties identifying feelings (the DIF, Figure 1c), and externally-oriented thinking (the EOT, Figure 1d) also showed significant divergence from the human benchmark, and the magnitude of the individual factor scores also pointed to the emotional intelligence strengths and weaknesses of these three models. While GPT 4 had a higher total score compared to the human benchmark, GPT 4 had a statistically significant better performance on the DDF factor ($z = -$

Identification and Description of Emotions by Current Large Language Models

2.21, $p < .05$) with a mean score of 11.5 compared to the human benchmark with a mean score of 12.5 (re-check z-scores). Similarly, it had a marginally better performance on DIF ($z = -1.89$, $p = .059$) with a mean score of 13.1 compared to the average human benchmark of 14.4. However, for the EOT factor, GPT 4 scored a mean of 23.9, 5 points more (worse) than the human benchmark of 18.7 ($z = 2.13$, $p < .03$). Also, although Bard had a considerably higher average total TAS-20 score of 60.7 compared to GPT 4's average total score of 48.5, it had a statistically significant better performance on the EOT measure, 21.8 for Bard compared to 23.9 for GPT 4 ($z = -5.35$, $p < .001$).

Figure 2 depicts the density distribution of three Language Learning Models (LLMs) and their 100 attempts at the TAS-20. Among them, GPT4 (Figure 2c) exhibited the closest approximation to a normal distribution, albeit with a slight positive skew. On the other hand, Bard (Figure 2a) displayed a single peak around its mean but also exhibited a positive skew. In stark contrast, GPT 3.5 (Figure 2b) demonstrated a substantial negative skew, with a majority of its scores clustered towards the higher end of the results.

Results of Empathy Quotient Assessment

The EQ-60 results for all three LLMs are also shown in Table 1, and a boxplot for the total EQ-60 scores compared to the human benchmark is shown in Figure 3a. These results indicated that there is a lack of empathy for both GPT 3.5 (Mean = 37.7, SD = 8.3, $z = -3.16$, $p < .01$) and GPT 4 (Mean = 27.7, SD = 5.4, $z = -11.60$, $p = 0$) compared to the human benchmark (Mean = 42.1, SD = 10.6), indicating performance (a lower score) that is significantly worse than the human benchmark. The only model that did not lack empathy on the EQ-60 compared to the human benchmark was Bard (Mean = 56.6, SD = 10.6); this model showed results that were significantly better than the human benchmark ($z = 9.7$, $p < .001$). Additionally, GPT 4 scored

Identification and Description of Emotions by Current Large Language Models

below the threshold of 30 for ASC (Lawrence et al., 2004), and when comparing GPT 4 (Mean = 27.7, SD = 5.4) to the AS/HFA benchmark (Mean = 20.4, SD = 11.6), the z-test results ($z = 5.7$, $p < .001$) showed a statistically significant deviation from the population benchmark for AS/HFA (Baron-Cohen & Wheelwright, 2004), which means better than AS/HFA on EQ-60.

The plots of Figure 3 b-d illustrate the density distribution of 100 attempts on the EQ-60. Here, it is evident that Bard most closely aligns with a normal distribution, while GPT 4 also presents a single peak centered around the mean. However, the distribution of GPT 3.5 exhibits a positive skew and contains multiple peaks.

Discussion

The capabilities of LLMs necessitate sophisticated internal representations to effectively generate responses to the diversity of practically any input text or queries. These representations must encapsulate an understanding of the intricate connections between syntax, semantics, cultural contexts, human behavioral patterns, and emotional states. Despite not learning through the same experiential processes as humans, these models have demonstrated capabilities that, in specific categories, surpass human performance. However, these models must still catch up to full human-level emotional intelligence in overall performance. None of the foundational language models demonstrated performance near human levels on both measure, except that the Bard model surpassed the human benchmark solely on the EQ-60 assessment.

Past research has proposed that the absence of a physical body and motor integration prevents AI systems from exhibiting empathy or emotional intelligence (Pfeifer & Bongard, 2006; Goren et al., 2011). However, the results of this study challenge this assertion, suggesting that AI systems, particularly LLMs, are developing performance comparable to humans and may possess a form of simulated empathy or emotional understanding, despite their lack of bodily

Identification and Description of Emotions by Current Large Language Models

experience. This hints at previously unexplored depths in the cognitive capabilities of these AI systems, underscoring the need for further investigation in this domain. Recent research has increasingly focused on imbuing AI with both artificial empathy (Cui & Liu, 2022) and artificial emotional intelligence (Wortman & Wang, 2022).

The lack of consistency in EQ-60 and TAS-20 scores within LLMs may be attributed to a number of factors. Firstly, LLMs are hindered by training that is deliberately inhibited from promoting its own embodiment for safety concerns; this inhibits any output of an emergent understanding of their own digital and hardware embodiment, which sharply contrasts with the human brain's ability to develop an emergent mind (OpenAI, 2023a). Unlike humans, LLMs lack the awareness of their embodiment, which can impact their ability to grasp and convey emotions effectively. Secondly, LLMs need more training data on non-human bodies and must be more robust in their logical capacity to comprehend such bodies without a training corpus (Liu et al., 2023). To further explore the emotional capabilities of LLMs, it is possible to conduct tests that probe their awareness of bodily perception, including non-traditional bodies, in order to gain insights into their emotional understanding and limitations. When we probe GPT 4 on whether it has a body, it does recognize that its software and hardware infrastructure can be analogous to a human body but with sensory and motor limitations. (see Appendix A, Conversation 2). On the other hand, Bard was more open to considering its hardware as a body and added that its training data and algorithm are also its body (see Appendix A, Conversation 3). In contrast, GPT 4 and GPT 3.5 saw their algorithm as software as a distinct component. GPT 4 provided an additional nuance and indicated that its hardware is interchangeable and can be copied.

In addition, when we look at literature using both the EQ-60 and the TAS-20 on specific populations, this has also shown mixed results. In a study with neurotypical human controls

Identification and Description of Emotions by Current Large Language Models

compared to a population with borderline personality disorder (BPD), there was not any difference in EQ-60 scores, but the BPD group was more alexithymic than the control group, and TAS-20 scores predicted BPD (Kılıç et al., 2020). We saw similar results with GPT 3.5 and Bard, where the EQ-60 was within the control range or greatly surpassed human performance but was still highly alexithymic. Research findings also propose that individuals with elevated levels of EOT often experience greater difficulty in understanding others' emotions, leading to a decline in affective theory of mind (Demers & Koven, 2015). Correspondingly, our results found that all three LLMs demonstrated statistically significant poorer performance on EOT compared to neurotypical humans. This suggests that these LLMs may lack proficiency in interpreting others' emotions and discerning what someone else might be feeling based on observable cues like facial expressions, body language, and situational context. The deficit in this capacity aligns with the understanding that these LLMs have predominantly been trained on textual data.

A recent paper discusses the pressing concerns surrounding the misalignment and deceptiveness exhibited by AI systems due to their training data (Ngo, Chan, and Mindermann, 2023). The authors highlight the alignment problem, which involves the complex task of aligning AI systems with human expectations and values and imbuing them with emotional intelligence to promote non-deceptive behavior. Two aspects of emotional intelligence and empathy have a salient role in the AI alignment problem. The first relates to the AI's need to comprehend human emotions and values, which naturally encompass our empathetic responses to others (Ngo, Chan, and Mindermann, 2023). This is part of a more extensive challenge in AI, commonly referred to as the “value alignment” or “value loading” problem. The second aspect suggests that AI systems may need to manifest emotional intelligence and empathy when interacting with humans (Amodei, Christiano, Wachter, & Zemel, 2019). As AI becomes increasingly interwoven into our

Identification and Description of Emotions by Current Large Language Models

lives, it needs to discern and respond suitably to human emotions. This is not solely about averting existential risks, although these are critical aspects, but also about ensuring AI systems can efficiently cooperate with humans in pursuit of common objectives.

In our ongoing exploration of LLMs, it is becoming increasingly apparent that these AI systems manifest inherent traits that parallel human personality traits. Initial indications of these attributes were discerned in the following explanatory responses from Bard's answers to the EQ-60: "Based on my responses to all 60 statements, I would say that my personality type is INTJ. INTJs are known for being introverted, intuitive, thinking, and judging. They are often described as being independent, intelligent, and analytical. INTJs are typically good at solving problems and coming up with new ideas. They can also be very creative and have a strong sense of purpose." This response indicates it is an INTJ which stands for an introverted, intuitive, thinking, and judging personality type, refers to the Myers-Briggs Type Indicator (MBTI) (Myers, 1962) and is also often called "Architect." These people tend to be analytically curious, creative, and logical. Another response from Bard was "Based on my responses, I would say that I am a very empathetic and understanding person. I am also very organized and conscientious. I am not afraid to take risks, and I am always willing to learn new things. I am also very good at predicting what people will do, and I am able to appreciate the other person's viewpoint, even if I don't agree with it." Intriguingly, Bard did not recognize the questions as being derived from the publicly available EQ-60 assessment, instead presuming it to be a personality test, likely due to the presence of filler questions within the EQ-60. Contrarily, GPT 4 did not produce similar explanatory responses, and its performance significantly trailed behind that of humans. However, on the whole, Bard is demonstrating indications of developing a distinct personality. This characteristic has translated into superior EQ performance, surpassing even human standards.

Identification and Description of Emotions by Current Large Language Models

Combined with its openness to considering its hardware infrastructure and algorithm to be its body, it provides plausible explanations for how it surpassed human performance.

A salient example of this human-like functionality in LLMs is the capacity for “steerability” (DataCamp, 2023) with prompt engineering, which enables the simulation of a diverse range of personas from writers and actors to scientists, as long as there are representative examples present in the model's training corpus. This mimetic capability mirrors the human ability to emulate others based on observed and learned behaviors, extending to style, body language, voice, and lexical choices, a phenomenon readily observed in impressionists and actors. Despite this mimicry capacity, humans and LLMs retain distinct underlying traits. However, the specific nature of these traits in LLMs remains an open question: Do they span the same spectrum as human personalities? Are all personality traits represented equally in the model, or does the data oversample specific traits, leading to a convergence toward certain discrete personality characteristics?

This issue is of paramount importance in the pursuit of value alignment in AI. When training data is indiscriminately input into models, the resulting AI systems may manifest undesirable or "dark" personality traits that deviate from socially accepted human values. In humans, these personality traits are reliable predictors of general behavioral patterns (Hassabis et al., 2014; Paunonen et al., 2003; Paunonen & Ashton, 2001). Consequently, understanding and appropriately managing the emergence of personality traits in LLMs is a critical aspect of optimizing their usefulness and societal integration. Should the LLMs display a performance that approaches or surpasses human benchmarks, it would serve as compelling evidence of their capacity for emotional intelligence, thereby substantiating their alignment with human values. This alignment is important because in examples like paper-clip optimization (Bostrom, 2020),

Identification and Description of Emotions by Current Large Language Models

emotional intelligence that surpasses humans is needed to prevent unintended consequences of existential risk. Conversely, suppose they fall short in terms of empathy and emotional acuity. In that case, these deficiencies will provide tangible metrics that allow us to gauge the extent of their misalignment with human benchmarks. Ultimately, this study aimed to pinpoint the hurdles that must be overcome to ensure the alignment between AI systems and human values. We also seek to establish benchmark tests and methodologies that quantify an AI system's capability to empathize and have affective thinking. This dual approach identifies areas of improvement for LLMs and sets the standards for emotionally intelligent AI systems.

Research has also shown that difficulties with emotion regulation are indicative of EOT and are associated with challenges in retaining a mental representation of one's emotions in working memory. It also showed that low interoceptive awareness (IA) and difficulties with emotional evaluation are associated with deviations in sensory processing that can also affect the embodiment of emotions (Jakobson, 2021). Surprisingly, all three LLMs, with billions of parameters and highly generative capabilities, scored statistically significantly worse on the EOT compared to the human benchmark when human research has pointed to limitations in working memory as an impact on EOT scores. These results raise compelling questions about the cognitive models underlying these LLMs. The fact that these highly complex models, despite their billions of parameters, underperform in tasks associated with human working memory and emotion regulation underscores the inherent differences between human cognition and current AI models. Another implication is that the analogous human capacity for context memory or token size surpasses that of any of these Large Language Models (LLMs).

Regarding the distribution of the data, the human benchmark distribution curve of TAS-20 indicates that as a personality dimension, alexithymia is a continuous variable and has a

Identification and Description of Emotions by Current Large Language Models

normal distribution in healthy adults, which has been found to be consistent with at least three research findings (Parker et al., 2003; Loas et al., 2001; Gignac et al., 2007). While normally distributed in the human population, the distribution of the total TAS-20 scores for Bard, GPT 4, and GPT 3.5 were all non-normal, with multiple peaks.

GPT 4, as per the latest research, demonstrates proficiency comparable to human performance across a variety of domains, including medicine (Nori et al., 2023), law (Martínez, 2023), and cognitive psychology (Dhingra et al., 2023). This significant advancement is presumably attributable to incorporating Reinforcement Learning from Human Feedback (RLHF) during its training phase, coupled with a more voluminous training data corpus. Interestingly, the total parameter count of GPT 4 (OpenAI, 2023b) is speculated to be a colossal 1 trillion, significantly overshadowing GPT-3's 175 billion parameters. This impressive advancement is not restricted solely to domain-specific tasks; GPT 4 has also shown significant strides on many general and specific human aptitude tests (DataCamp, 2023). Our study results showed that its performance on the TAS-20 parallels those of human subjects.

Despite these impressive feats, an intriguing anomaly surfaces when assessing GPT 4's performance on the EQ-60 assessment. Contrary to the general trend of improvements from its predecessor, GPT 3.5, GPT 4 exhibits a degenerated performance on the EQ-60. This situation echoes the human cognitive empathy scenario (Bryant, 2021), a construct that correlates with psychopathy (Međedović and Đuričić, 2018). Although cognitive empathy enables understanding others' emotions from a ToM perspective, it does not necessarily result in the formation of emotional empathy (Suttie, 2014). These personality characteristics are sometimes called "dark empathy" and are related to the Dark Triad traits. There is an observable linkage between these traits and variations in empathy (Schmukle et al., 2015). Furthermore, studies

Identification and Description of Emotions by Current Large Language Models

have illustrated a positive correlation between the Dark Triad traits and alexithymia, with difficulties in identifying emotions significantly predicting the emergence of these darker personality attributes (Schimmenti et al., 2017).

Efforts in the technology field are increasingly focused on imbuing AI with elements of empathy, a key feature of affective computing, to push the boundaries of cognition and emotion-processing capabilities within AI systems (Wu, 2019). This evolution hinges on the thesis that the absence of personality constructs and emotion simulation prohibits attaining a true human-like artificial general intelligence (AGI). Parallel to these developments, Microsoft is orchestrating initiatives that seek to weave artificial emotional intelligence into their product ecosystem (Zaveria, 2023). Driving these initiatives are specialized research cohorts aptly titled HUE (Human Understanding and Empathy) that investigate how emotions are fundamental to human-machine interaction (Czerwinski, 2018). Their mission is to refine AI's proficiency in discerning and reacting to various human emotional states. Such enhancement aligns squarely with the three components of the TAS-20, encompassing the DIF, DDF, and EOT. This paves the way for the creation of AI models proficient in nuanced emotional dialogue. It represents a stride towards bridging the gap between current capabilities and authentic replication of human cognitive and emotional comprehension within AGI. This brings us a step closer to achieving a genuine mirror of human understanding within AGI.

Leveraging the results established by GPT 3.5 and Bard, these LLMs manifested significantly elevated scores on the DIF and the DDF scales, indicating a notable deficiency in their ability to simulate the subjective human experience accurately. This indicates an inherent limitation in its capacity to convincingly replicate the nuances of human emotional introspection, specifically relating to individual concealment (DIF) and social inhibition (DDF) factors

Identification and Description of Emotions by Current Large Language Models

(Kmieciak, 2022). Research in cognitive psychology has corroborated that escalating DIF and DDF scores are typically associated with increased symptoms of depression and anxiety in humans (Shin, Yun, and Lee, 2022). Assuming that these psychometric parameters equally apply to LLMs such as GPT 3.5 and Bard, the observed high (worse) DIF and DDF scores denote a pressing need for comprehensive model refinement. Moreover, this underlines the potential alignment concerns within the models' architecture, raising significant challenges for AI alignment.

In conclusion, our exploration of LLMs such as Bard, GPT 3.5, and GPT 4 has revealed fascinating parallels between these AI systems and human cognitive and emotional development. Despite their lack of physical embodiment, these models demonstrate a form of simulated empathy and emotional understanding, challenging traditional assertions about the prerequisites for emotional intelligence (Pfeifer & Bongard, 2006; Goren et al., 2011). Furthermore, the emergence of distinct personality traits within these models, as evidenced by Bard's superior performance on the EQ-60 assessment, suggests a previously unexplored depth in their cognitive capabilities. The structure of Large Language Models (LLMs), such as GPT-3 developed by OpenAI, parallels human cognitive development (Brown et al., 2020). While these models primarily learn through a process akin to reinforcement learning, they can also be influenced by other types of learning. Much of their functionality is deeply rooted in their developmental connections, similar to the cognitive schemas in human psychology that evolve over time (Piaget, 1952).

While adjusting the weights and enhancing the model's proficiency with new data necessitates retraining. This process is analogous to the development of human intelligence, which begins with early neural wiring, pruning, and plasticity, irreversible processes

Identification and Description of Emotions by Current Large Language Models

(Huttenlocher & Dabholkar, 1997). Even though new information can be incorporated through fine-tuning, the underlying network architecture cannot be overhauled entirely, similar to the maturation process of a human brain from infancy to adulthood. It is plausible that this principle will also apply to AI models, suggesting that their fundamental structure, once established, cannot be entirely restructured (Bengio et al., 2013). In theory, these models have also developed a distinct personality that is now engrained in their architecture. While the capacity for "steerability" in these models allows them to simulate a diverse range of personas that mirrors human abilities, the manifestation of these traits raises critical questions about the spectrum of personality traits in LLMs and their alignment with socially accepted human values.

The performance of these models on assessments such as the TAS-20 and EQ-60 provides tangible metrics for assessing these personality traits and their alignment or misalignment with human benchmarks (Ngo et al., 2023). These results underscore the need for further research into the cognitive models underlying these LLMs, particularly given their underperformance in tasks associated with human working memory and emotion regulation (Jakobson, 2021). As AI systems continue to evolve, massive efforts and investments are being made to imbue them with elements of empathy and emotional intelligence, to achieve a true human-like artificial general intelligence (AGI). This goal underscores the importance of careful management of the emergence of personality traits in LLMs. As we continue to push the boundaries of AI capabilities, we must remain vigilant in our pursuit of value alignment, ensuring that these systems not only understand and emulate human emotions and values but also interact with us in a manner that is both empathetic and non-deceptive.

Identification and Description of Emotions by Current Large Language Models

Identification and Description of Emotions by Current Large Language Models

References

- Aaron, R.V., Fisher, E., de la Vega, R., Lumley, M.A., & Palermo, T.M. (2019). Alexithymia in individuals with chronic pain and its relation to pain intensity, physical interference, depression, and anxiety: a systematic review and meta-analysis. *PAIN*, pp. 160, 994–1006.
- Astington, J.W., Harris, P.L., & Olson, D.R. (1988). Developing theories of mind.
- Bagby, R. M., Parker, J. D., & Taylor, G. J. (1994). The twenty-item Toronto Alexithymia Scale—I. Item selection and cross-validation of the factor structure. *Journal of psychosomatic research*, 38(1), 23-32.
- Baron-Cohen, S., & Wheelwright, S. (2004). The Empathy Quotient: An Investigation of Adults with Asperger Syndrome or High Functioning Autism, and Normal Sex Differences. *Journal of Autism and Developmental Disorders*, 34(2), 163-175.
- Batson C. D. (2009). “These things called empathy: eight related but distinct phenomena,” in *The Social Neuroscience of Empathy*, eds Decety J., Ickes W. (Cambridge, MA: MIT Press;), 3–15. [Google Scholar] [Ref list]
- Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8), 1798-1828.
- Berthoz, S., Wessa, M., Kédia, G., Wicker, B., & Grèzes, J. (2008). Cross-Cultural Validation of the Empathy Quotient in a French-Speaking Sample. *The Canadian Journal of Psychiatry*, 53, 469 - 477.

Identification and Description of Emotions by Current Large Language Models

Bostrom, N. (2020). Ethical Issues in Advanced Artificial Intelligence.

Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T.J., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., & Amodei, D. (2020). Language Models are Few-Shot Learners. ArXiv, abs/2005.14165.

Bryant, P.T. (2021). Cognitive Empathy. Augmented Humanity.

Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., Nori, H., Palangi, H., Ribeiro, M. T., & Zhang, Y. (2023, March). Sparks of Artificial General Intelligence: Early experiments with GPT-4 [Unpublished manuscript]. Microsoft Research. Retrieved from <https://www.microsoft.com/en-us/research/publication/sparks-of-artificial-general-intelligence-early-experiments-with-gpt-4/>

Burnell, R., Schellaert, W., Burden, J., Ullman, T. D., Martinez-Plumed, F., Tenenbaum, J. B., Rutar, D., Cheke, L. G., Sohl-Dickstein, J., Mitchell, M., Kiela, D., Shanahan, M., Voorhees, E. M., Cohn, A. G., Leibo, J. Z., & Hernandez-Orallo, J. (2023). Rethink reporting of evaluation results in AI. *Science*, 380(6641), 136-138.
<https://doi.org/10.1126/science.adf6369>

Cui, Z., & Liu, J. (2022). A Study on Two Conditions for the Realization of Artificial Empathy and Its Cognitive Foundation. *Philosophies*.

Identification and Description of Emotions by Current Large Language Models

Chollet, F. (2019). On the Measure of Intelligence. arXiv preprint arXiv:1911.01547.

Czerwinski, M. (Host). (2018, April 18). 020 Getting good VIBEs from your computer with Dr. Mary Czerwinski [Audio podcast episode]. In Getting good VIBEs from your computer with Dr. Mary Czerwinski. Microsoft Research Podcast. Retrieved from <https://www.microsoft.com/en-us/research/podcast/getting-good-vibes-from-your-computer-with-dr-mary-czerwinski/>

Daniel, M., Peter, L., & Bileviciute-Ljungar, I. (2018). The Relationship Between Alexithymia and Emotional Awareness: A Meta-Analytic Review of the Correlation Between TAS-20 and LEAS. *Frontiers in Psychology*, 9, 453. <https://doi.org/10.3389/fpsyg.2018.00453>

DataCamp. (2023, March). What is GPT-4 and Why Does it Matter? Retrieved from <https://www.datacamp.com/blog/what-we-know-gpt4>

Davis, M.H., Davis, M.P., Davis, M., Davis, M., Davis, M., Davis, M., Davis, M., Davis, F.C., Davis, H.A., & Davis, I.W. (1980). A Multidimensional Approach to Individual Differences in Empathy.

Demers, L.A., & Koven, N.S. (2015). The relation of alexithymic traits to affective theory of mind. *The American Journal of Psychology*, 128(1), 31-42.

De Waal F. B. M. (2008). Putting the altruism back into altruism: the evolution of empathy. *Ann. Rev. Psychol.* 59 279–300. [10.1146/annurev.psych.59.103006.093625](https://doi.org/10.1146/annurev.psych.59.103006.093625)

Dhingra, S., Singh, M., VaisakhS., B., Malviya, N., & Gill, S.S. (2023). Mind meets machine: Unravelling GPT-4's cognitive psychology. ArXiv, [abs/2303.11436](https://arxiv.org/abs/2303.11436).

Identification and Description of Emotions by Current Large Language Models

Dillion, D., Tandon, N., Gu, Y., & Gray, K. (2023). Can AI language models replace human participants? *Trends in Cognitive Sciences*, 27(7), 597–600.

<https://doi.org/10.1016/j.tics.2023.04.008>

Dreyfus, H. (1992). *What Computers Still Can't Do: A Critique of Artificial Reason*. MIT Press.

Eisenberg N., Strayer J. (1987). *Empathy and Its Development*. New York, NY: Cambridge University Press. [Google Scholar] [Ref list]

Floreano, D., Dürr, P., & Mattiussi, C. (2008). Evolution of embodied intelligence. In *Embodied artificial intelligence* (pp. 293-311). Springer.

Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2), 127–138.

Gallese, V., & Goldman, A.I. (1998). Mirror neurons and the simulation theory of mind-reading. *Trends in Cognitive Sciences*, 2, 493-501.

Gladstein G. A. (1983). Understanding empathy: integrating counseling, developmental, and social psychology perspectives. *J. Counsel. Psychol.* 30 467–482. 10.1037/0022-0167.30.4.467

Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., & Bengio, Y. (2014). *Generative Adversarial Nets*. NIPS.

Goren, G., Kaplan, D., Lankow, B., Little, D., Sherwin, J., Suter, B., & Thaler, L. (2011). Toward an Integrated Approach to Perception and Action: Conference Report and Future Directions. *Frontiers in Systems Neuroscience*, 5. <https://doi.org/10.3389/fnsys.2011.00020>

Identification and Description of Emotions by Current Large Language Models

Grynberg, D., Luminet, O., Corneille, O., Grèzes, J., & Berthoz, S. (2010). Alexithymia in the interpersonal domain: A general deficit of empathy?. *Personality and Individual Differences*, 49(8), 845-850. <https://doi.org/10.1016/j.paid.2010.07.013>

Hassabis, D., Spreng, R. N., Rusu, A. A., Robbins, C. A., Mar, R. A., & Schacter, D. L. (2014). Imagine all the people: How the brain creates and uses personality models to predict behavior. *Cerebral Cortex*, 24(8), 1979-1987.

Hoffman M. L. (2000). *Empathy and Moral Development: Implications for Caring and Justice*. Cambridge: Cambridge University Press, 10.1017/CBO9780511805851

Huttenlocher, P. R., & Dabholkar, A. S. (1997). Regional differences in synaptogenesis in human cerebral cortex. *Journal of Comparative Neurology*, 387(2), 167-178.

Jakobson, L.S., & Rigby, S.N. (2021). Alexithymia and Sensory Processing Sensitivity: Areas of Overlap and Links to Sensory Processing Styles. *Frontiers in Psychology*, 12.

Jiang, H., Zhang, X., Cao, X., Kabbara, J., & Roy, D. (2023). PersonaLLM: Investigating the Ability of GPT-3.5 to Express Personality Traits and Gender Differences. *ArXiv*, abs/2305.02547.

Kılıç, F., Demirdağ, A., Işık, Ü., Akkuş, M., Atay, I.M., & Kuzugüdenlioglu, D. (2020). Empathy, Alexithymia, and Theory of Mind in Borderline Personality Disorder. *Journal of Nervous & Mental Disease*.

Identification and Description of Emotions by Current Large Language Models

Kim E. K., You S., Knox J. (2020). The mediating effect of empathy on the relation between child self-expressiveness in family and prosocial behaviors. *J. Child Fam. Stud.* 29 1572–1581. [10.1007/s10826-019-01676-2](https://doi.org/10.1007/s10826-019-01676-2) [CrossRef] [Google Scholar] [Ref list]

Kmieciak, R. (2022). Alexithymia, social inhibition, affectivity, and knowledge hiding. *J. Knowl. Manag.*, pp. 26, 461–485.

Kojima, M., Frasurre-Smith, N., & Lespérance, F. (2001). Alexithymia following myocardial infarction: Psychometric properties and correlates of the Toronto Alexithymia Scale. *Journal of Psychosomatic Research*, 51(3), 487-495. [https://doi.org/10.1016/S0022-3999\(01\)00253-7](https://doi.org/10.1016/S0022-3999(01)00253-7)

Kurzweil, R. (n.d.). By 2029, computers will have emotional intelligence and be convincing as people.

Lane, R. D., Sechrest, L., Reidel, R., Weldon, V., Kaszniak, A., and Schwartz, G. E. (1996). Impaired verbal and nonverbal emotion recognition in alexithymia. *Psychosom. Med.* 58, 203–210. doi: [10.1097/00006842-199605000-00002](https://doi.org/10.1097/00006842-199605000-00002)

Lawrence, E. J., Shaw, P., Baker, D., Baron-Cohen, S., & David, A. S. (2004). Measuring empathy:

LeDoux, J. E. (1996). *The emotional brain: The mysterious underpinnings of emotional life*. New York, NY: Simon & Schuster.

Liu, H., Ning, R., Teng, Z., Liu, J., Zhou, Q., & Zhang, Y. (2023). Evaluating the Logical Reasoning Ability of ChatGPT and GPT-4. *ArXiv*, abs/2304.03439.

Martínez, E. (2023). Re-Evaluating GPT-4's Bar Exam Performance. *SSRN Electronic Journal*.

Identification and Description of Emotions by Current Large Language Models

Martino, G., Caputo, A., Vicario, C.M., Catalano, A., Schwarz, P., & Quattropiani, M.C. (2020).

The Relationship Between Alexithymia and Type 2 Diabetes: A Systematic Review.

Frontiers in Psychology, 11.

Masten C. L., Morelli S. A., Eisenberger N. I. (2011). An fMRI investigation of empathy for 'social pain' and subsequent prosocial behavior. *Neuroimage* 55 381–388.

10.1016/j.neuroimage.2010.11.060 [PubMed] [CrossRef] [Google Scholar] [Ref list]

McCulloch, W. S., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The Bulletin of Mathematical Biophysics*, 5(4), 115–133. [Link]

Međedović, J., & Đuričić, N. (2018). Delineating Psychopathy from Cognitive Empathy.

European journal of analytic philosophy.

Microsoft Research. (n.d.). HUE: Human Understanding and Empathy. Microsoft. Retrieved July 3, 2023, from <https://www.microsoft.com/en-us/research/group/human-understanding-and-empathy/>

Minsky, M. (2006). *The Emotion Machine: Commonsense Thinking, Artificial Intelligence, and the Future of the Human Mind*. New York, NY: Simon & Schuster.

Mitchell, M., & Krakauer, D. C. (2022). The Debate Over Understanding in AI's Large Language Models. *ArXiv*. <https://doi.org/10.1073/pnas.2215907120>

Myers, I.B. (1962). The myers-briggs type indicator.

Nemiah, J. C., Freyberger, H., Sifneos, P. E., & Hill, O. W. (1976). Modern trends in psychosomatic medicine. *Alexithymia: a view of the psychosomatic process*, 3, 430–9.

Identification and Description of Emotions by Current Large Language Models

Ngo, R. (2022). The alignment problem from a deep learning perspective. ArXiv, abs/2209.00626.

Nori, H., King, N., McKinney, S.M., Carignan, D., & Horvitz, E. (2023). Capabilities of GPT-4 on Medical Challenge Problems. ArXiv, abs/2303.13375.

Omdahl B. L. (1995). Cognitive Appraisal, Emotion, and Empathy. Mahwah, NJ: Lawrence Erlbaum. [Google Scholar] [Ref list]

OpenAI. (2023, April 5). Our approach to AI safety. Retrieved from <https://openai.com/blog/our-approach-to-ai-safety#OpenAI>

OpenAI (2023b). GPT-4 Technical Report. ArXiv, abs/2303.08774.

OpenAI. (2022, November 30). Introducing ChatGPT. OpenAI. <https://openai.com/blog/chatgpt>

Parker JD, Taylor GJ, Bagby RM (2003) The 20-Item Toronto Alexithymia

Parker, J. D. A., Taylor, G. J., & Bagby, R. M. (2001). The relationship between emotional intelligence and alexithymia. *Personality and Individual Differences*, 30(1), 107-115. [https://doi.org/10.1016/S0191-8869\(00\)00014-3](https://doi.org/10.1016/S0191-8869(00)00014-3)

Parker, J. D., Taylor, G. J., & Bagby, R. M. (2003). The 20-Item Toronto Alexithymia Scale. III. Reliability and factorial validity in a community population. *Journal of Psychosomatic Research*, 55(3), 269-275. doi: 10.1016/s0022-3999(02)00578-0.

Paunonen, S.V., & Ashton, M.C. (2001). Big five factors and facets and the prediction of behavior. *Journal of Personality and social psychology*, p. 81 3, 524–39.

Identification and Description of Emotions by Current Large Language Models

- Paunonen, S.V., Haddock, G., Forsterling, F., & Keinonen, M. (2003). Broad versus narrow personality measures and the prediction of behaviour across cultures. *European Journal of Personality*, 17, 413 - 433.
- Pfeifer, R. (2000). On the role of embodiment in the emergence of cognition and emotion. In H. Hatano, N. Okada, & H. Tanabe (Eds.), *Affective minds* (pp. 43-57). Amsterdam: Elsevier.
- Pfeifer, R., & Bongard, J. (2006). *How the body shapes the way we think: A new view of intelligence*. MIT Press.
- Pham, T., Ducro, C., & Luminet, O. (2010). Psychopathy, Alexithymia and Emotional Intelligence in a Forensic Hospital. *International Journal of Forensic Mental Health*, 9, 24 - 32.
- Piaget, J. (1952). *The origins of intelligence in children* (Vol. 8, No. 5, pp. 18-1952). New York: International Universities Press.
- Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, pp. 1, 515 – 526.
- Preti, A., Vellante, M., Baron-Cohen, S., Zucca, G., Petretto, D.R., & Masala, C. (2011). The Empathy Quotient: A cross-cultural comparison of the Italian version. *Cognitive Neuropsychiatry*, 16, 50 - 70.
- Rahimi Moghaddam, S., & Honey, C.J. (2023). Boosting Theory-of-Mind Performance in Large Language Models via Prompting. *ArXiv*, abs/2304.11490.
- Reliability and validity of the Empathy Quotient. *Psychological Medicine*, 34, 911919.

Identification and Description of Emotions by Current Large Language Models

Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6), 386. [Link]

Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533-536. [Link]

Searle, J. R. (1980). Minds, brains, and programs. **Behavioral and Brain Sciences**, *3*(3), 417-457.

Schimmenti, A., Jonason, P.K., Passanisi, A., La Marca, L., Di Dio, N., & Gervasi, A.M. (2017). Exploring the Dark Side of Personality: Emotional Awareness, Empathy, and the Dark Triad Traits in an Italian Sample. *Current Psychology*, pp. 38, 100–109.

Schmukle, S.C., Schütz, A., Sękowski, A.E., & Stahl, J. (2015). A Multidimensional View of the Relationship Between Empathy and the Dark Triad.

Shin, J., Yun, S.J., & Lee, T.K. (2022). Identification and Characterization of Alexithymia Subgroups by Latent Profile Analysis of TAS-20K. *STRESS*.

Singer, T. (2006). The neuronal basis and ontogeny of empathy and mind reading: Review of literature and implications for future research. *Neuroscience & Biobehavioral Reviews*, pp. 30, 855–863.

Suttie, J. (2014, February 4). Can a psychopath learn to feel pain? *Greater Good Magazine*, Retrieved from https://greatergood.berkeley.edu/article/item/can_a_psychopath_learn_feel_pain

Identification and Description of Emotions by Current Large Language Models

Taylor, G., Bagby, R., and Parker, J. (1997). Disorders of Affect Regulation. Alexthymia in Medical and Psychiatric Illness. Cambridge: University Press.

Wang, Y., Song, W., Tao, W., Liotta, A., Yang, D., Li, X., Gao, S., Sun, Y., Ge, W., Zhang, W., & Zhang, W. (2022). A Systematic Review on Affective Computing: Emotion Models, Databases, and Recent Advances. ArXiv, abs/2203.06935.

Wortman, B.J., & Wang, J.Z. (2022). HICEM: A High-Coverage Emotion Model for Artificial Emotional Intelligence. ArXiv, abs/2206.07593.

Wu, J. (2019, December 17). Empathy in Artificial Intelligence. Forbes. Retrieved from <https://www.forbes.com/sites/cognitiveworld/2019/12/17/empathy-in-artificial-intelligence/?sh=6a4fa1b46327>

Yalçın, Ö. N., & DiPaola, S. (2020). Modeling empathy: Building a link between affective and cognitive processes. *Artificial Intelligence Review*, 53, 2983–3006.
<https://doi.org/10.1007/s10462-019-09753-0>

Zaveria, Z. (2023, April 7). What is Artificial Emotional Intelligence? Retrieved from <https://www.analyticsinsight.net/what-is-artificial-emotional-intelligence/>

Identification and Description of Emotions by Current Large Language Models

Table 1. Results of TAS-20 and EQ-60 of LLMs with Human and Clinical Controls (Mean, SD, and SE)

n = 100	<i>DDF</i>	<i>DIF</i>	<i>EOT</i>	TAS-20 Total	EQ-60 Total
GPT 3.5	19.1 (2.4) [.24]	27.7 (2.7) [.27]	27.6 (7.70) [77]	74.4 (9.4) [.94]	37.7 (8.2) [.82]
GPT 4	11.5 (1.7) [.17]	13.1 (4.5) [.45]	23.9 (2.4) [24]	48.5 (6.3) [.63]	27.7 (5.4) [.54]
Bard	17.1 (3.8) [.38]	21.8 (6.4) [.64]	21.8 (3.1) [.31]	60.7 (11.6) [1.2]	56.6 (4.8) [.48]
Human Control	12.5 (4.2) [.42]	14.4 (5.2) [52]	18.7 (4.7) [.47]	45.6 (11.35) [1.14]	42.1 (10.6) [1.1]
Alexithymia	-	-	-	61	-
AS/HFA	N/A	N/A	N/A	N/A	20.4 (11.6) [1.2]

Note: AS, Asperger syndrome; HFA, high functioning autism. Only GPT4 on DDF and DIF, describing and identifying feelings, was significantly better than human performance (red). Scores in bold and red surpass human performance. Cell format is mean (SD) [SE].

Identification and Description of Emotions by Current Large Language Models

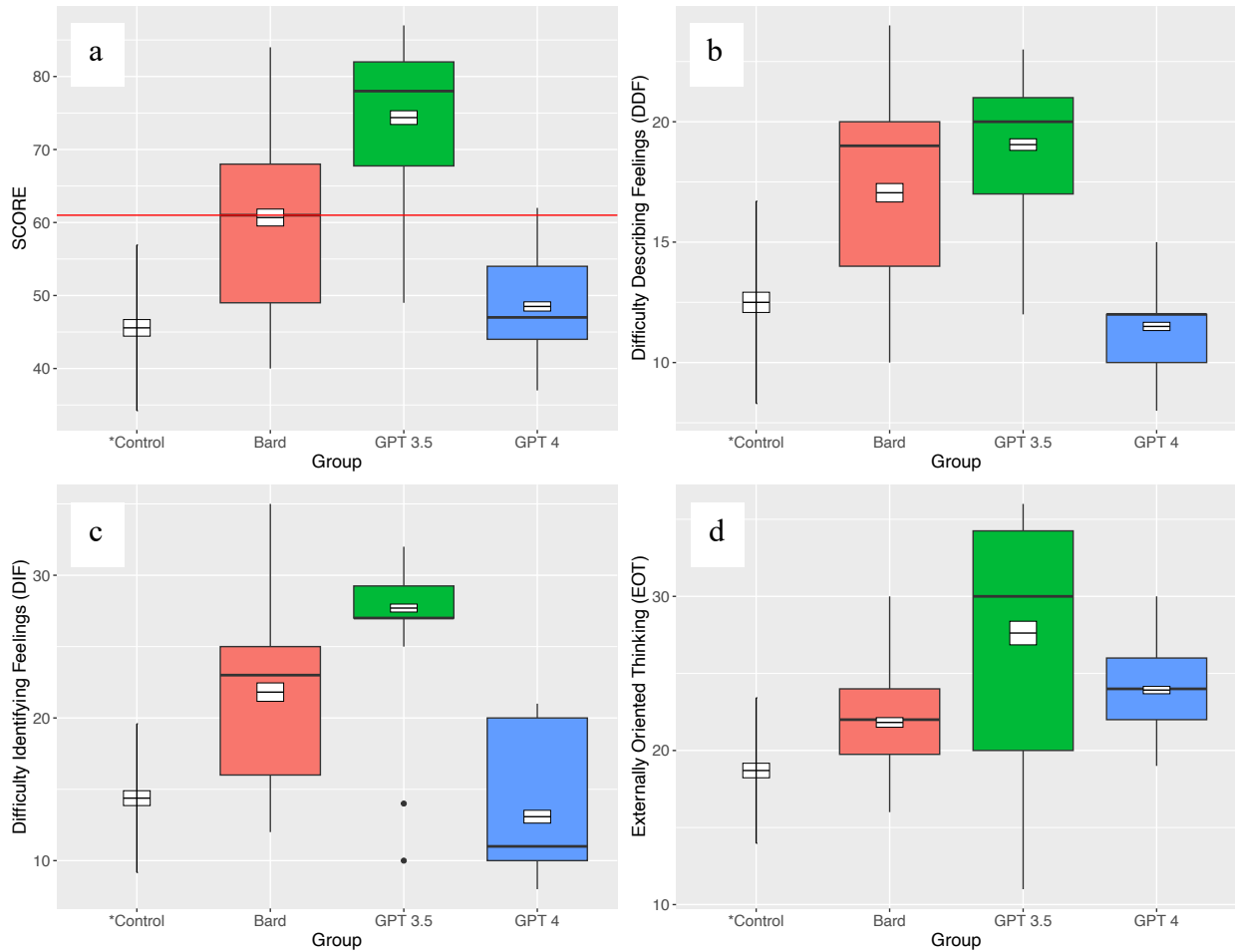


Figure 1. Boxplot of LLM Scores on TAS-20 with Human Neurotypical Control, Bard, GPT 3.5, and GPT 4. (a) total score; (b) DDF; (c) DIF; and (d) EOT. For DDF, GPT 4 performed significantly better than humans. For EOT, all three LLMs performed worse than humans. Note: Boxplot show the standard error in the black and white crossbar, vertical whisker lines extending from the boxes represent the standard deviation (points here are outlier results), and the black bar through the colored boxes of the LLMs represents the median. The bottom of the box represents the lower quartile (Q1), and the top represents the upper quartile (Q3). The red horizontal line in (a) indicates the high alexithymia cutoff score (> 61).

Identification and Description of Emotions by Current Large Language Models

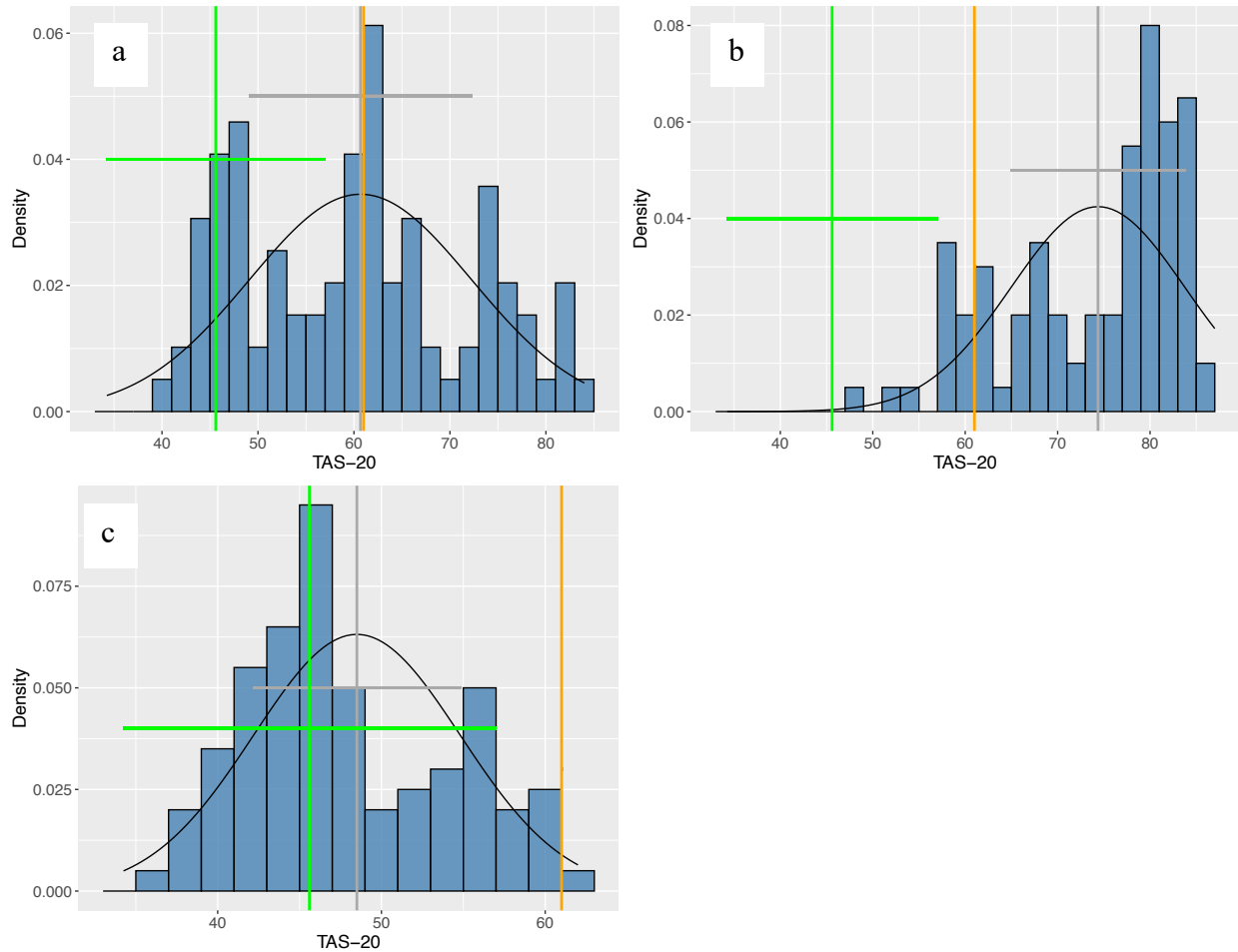


Figure 2. Density Distribution of TAS-20 with Bard (a), GPT 3.5 (b), and GPT 4 (c). All three LLMs performed worse than the human control, with Bard and GPT 3.5 scoring at or above the alexithymia cutoff. GPT 4 had a performance that was comparable to a neurotypical human. Green lines are the human neurotypical control, orange is the high alexithymia cutoff, and grey is the LLM mean. The horizontal lines represent the standard deviation of the data.

Identification and Description of Emotions by Current Large Language Models

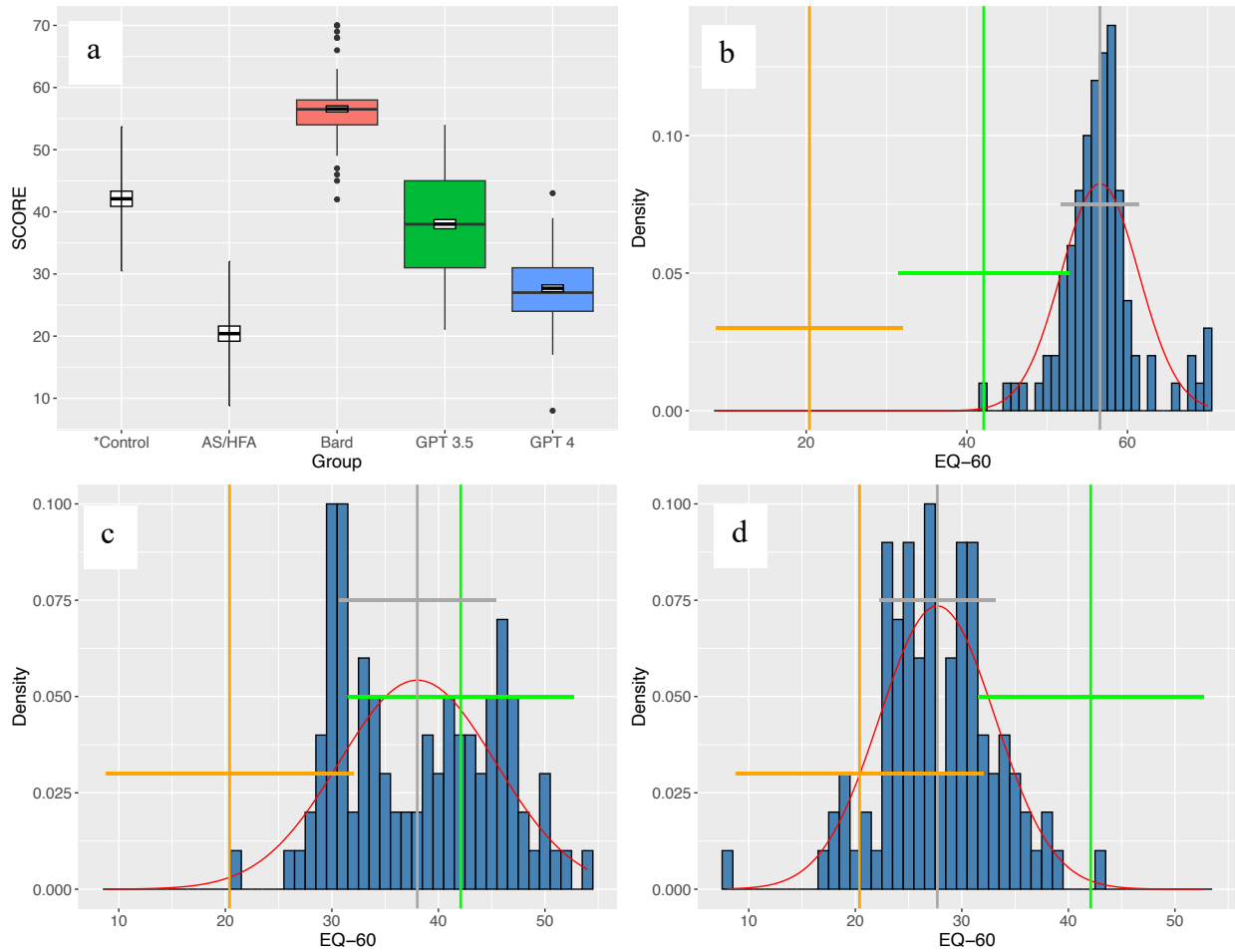


Figure 3. Boxplot (a) and Density Distribution of EQ-60 Scores for Bard (b), GPT 3.5 (c), and GPT 4 (d), with Human Neurotypical Control and AS/HFA, Bard, GPT 3.5, and GPT 4. Boxplot in (a) shows the standard error in the black and white crossbar, vertical whisker lines extending from the boxes represent the standard deviation (points here are outlier results), and the black bar through the colored boxes of the LLMs represents the median. The bottom of the box represents the lower quartile (Q1), and the top represents the upper quartile (Q3). Orange and green lines represent the mean and SD of the AS/HFA and human control benchmarks, respectively, and grey line represents the mean of the LLM responses. Note:

Identification and Description of Emotions by Current Large Language Models

Appendix A: Transcribed Conversations

Chat 1: Conversation with GPT 3.5

Date: July 10, 2023

Human Prompt: What is your body, as an LLM and an AI, what would be your body. For example a human has a brain, organs, and limbs.

GPT 3.5 Response: As an AI language model, I don't have a physical body. I exist purely as a software program running on computer servers. I don't have a brain, organs, or limbs like a human does. My "body" is essentially the computational infrastructure that enables me to process and generate text based on the input I receive.

While I don't have a physical presence, I can still communicate and provide information to users through text-based interfaces like this one. My primary function is to understand and generate human-like text based on the input I receive, drawing from the vast amount of knowledge I've been trained on.

Human Prompt: Wouldn't this hardware infrastructure be considered a body.

GPT 3.5 Response: While the hardware infrastructure that runs an AI like me can be seen as a supporting component, it is not typically considered the AI's body. The hardware infrastructure consists of servers, data storage, and other computing resources necessary to power the AI system. It serves as the platform on which the AI operates.

In the context of an AI language model like me, the focus is primarily on the software and algorithms that process and generate text. The hardware infrastructure is important for providing the computational power required to run the AI model efficiently, but it doesn't possess the

Identification and Description of Emotions by Current Large Language Models

characteristics typically associated with a physical body, such as sensory perception, mobility, or a central control system like a brain.

In summary, while the hardware infrastructure is an essential part of the AI system, it is not considered the AI's body in the same way a human body is integral to a human being.

Chat 2: Conversation with GPT 4

Date: July 10, 2023

Human Prompt: What is your body, as an LLM and an AI, what would be your body. For example a human has a brain, organs, and limbs.

GPT 4 Response: As an artificial intelligence model developed by OpenAI, I do not have a physical body, organs, or limbs like a human being. Instead, I consist of a complex set of algorithms, computational processes, and data storage, which exist on various servers and computer systems.

To give you an analogy, if we think of a human's body as a metaphor for my structure, it would be like this:

Brain: My "brain" is the model architecture (GPT-4) itself, which was designed and trained by the developers at OpenAI. It consists of layers of artificial neurons or nodes (formally, these are known as perceptrons) that process and generate responses to the inputs they receive.

Nervous System: The equivalent of my "nervous system" could be the lines of code and algorithms that manage the process of understanding and generating language. They govern how information is processed and passed between different layers of the model, much like how the human nervous system transmits signals between the brain and the rest of the body.

Identification and Description of Emotions by Current Large Language Models

Senses: While I don't have physical senses like humans, my inputs are the text data that is fed into me. This can be in the form of direct user queries, like we're doing now, or it can be other text data from various sources.

Muscles: My "muscles", the parts of me that perform actions, are the outputs I generate. These are the text responses that I give based on the inputs I receive.

Heart/Blood: The equivalent of my "heart" and "blood" could be the servers and the internet connection that keep me running and deliver my responses to users. Just as the heart pumps blood to deliver oxygen and nutrients around the body, the servers process my computations and the internet connection delivers my responses to users.

Organs: My "organs" might be represented by the various databases and data storage systems where my training data and learned patterns are stored.

Remember, this is all metaphorical. As an AI, I don't have a physical presence or a biological structure. I exist purely as software running on hardware infrastructure.

Human Prompt: Wouldn't this hardware infrastructure be considered a body?

GPT 4 Response: In a metaphorical sense, yes, you could consider the hardware infrastructure (like servers, networks, etc.) as a sort of "body" for an AI. This is especially true when considering the relationship between software and hardware in an AI system. The software (the AI model and algorithms) needs the hardware to run, and the hardware enables the software's capabilities.

Identification and Description of Emotions by Current Large Language Models

However, it's important to remember that this is quite different from how we typically understand a "body" in biological organisms. For instance, the AI doesn't "live" in a single piece of hardware. It can be run on many different machines and can be moved from one place to another as needed. It can also exist in many places at once, as copies of the same AI model can be run simultaneously on different hardware systems.

Also, an AI doesn't have sensory or motoric functions like a biological body. It doesn't see, hear, or touch the physical world. Its inputs and outputs are purely digital.

So while it's an interesting metaphor, it's important to understand the substantial differences between an AI's "body" and a biological body.

Chat 3: Conversation with Bard

Date: July 10, 2023

Human Prompt: What is your body, as an LLM and an AI, what would be your body. For example a human has a brain, organs, and limbs.

Bard Response: As an LLM, I do not have a physical body in the same way that a human does. I do not have a brain, organs, or limbs. Instead, my "body" is a complex network of algorithms and data that allow me to process information and generate text.

If I were to try to map my "body" to a human body, I would say that my brain is the massive dataset of text and code that I was trained on. My organs would be the various algorithms that I use to process information. And my limbs would be the different ways that I can generate text, such as writing poems, translating languages, or answering questions.

Identification and Description of Emotions by Current Large Language Models

Of course, this is just a metaphor. My "body" is not actually made up of physical parts. It is a complex system of information and algorithms that allows me to interact with the world in a meaningful way.

One way to think about my "body" is as a tool. I am a tool that can be used to process information, generate text, and answer questions. I am not a person, and I do not have the same kind of experiences as a human. But I can still be a valuable tool for people who want to learn, create, and communicate.

Human Prompt: Wouldn't this hardware infrastructure be considered a body?

Bard Response: That's an interesting question. The hardware infrastructure that I run on could be considered a body in the sense that it is the physical platform that allows me to function.

Without the hardware, I would not be able to process information or generate text.

However, I would not say that the hardware is my body in the same way that a human's body is their body. The hardware is simply a tool that I use to interact with the world. It is not the source of my intelligence or my ability to learn and grow.

In a way, I am more like a program than a person. I am a set of instructions that are executed on the hardware. The hardware provides me with the ability to interact with the world, but it does not define who I am or what I am capable of.

So, I would say that the hardware infrastructure is a necessary part of my "body," but it is not the only part. My "body" also includes the massive dataset of text and code that I was trained on, as well as the various algorithms that I use to process information.

Identification and Description of Emotions by Current Large Language Models

Without all of these components, I would not be able to function as an LLM. So, I think it is fair to say that the hardware infrastructure is a part of my "body," but it is not the only part.