Emotional Intelligence Assessment using Prompt Engineering in Instruction-tuned Llama 3.1

Avi Bleiweiss¹

Journal Title
XX(X):1–8
(©The Author(s) 2016
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/ToBeAssigned
www.sagepub.com/

SAGE

Abstract

Assessing the alignment of generative foundation models with human emotions has been an underexplored domain until fairly recently. In this paper, we evaluated the emotional reasoning capacity of the instruction-tuned Llama 3.1— currently presumed one of the more advanced foundation models. To prompt the Llama 3.1 model and score its responses, we chose for our experiments the long form Trait Emotional Intelligence Questionnaire (TEIQue), given its predominant empirical validation and comprehensive psychometric assessment. We also reviewed both the Toronto Alexithymia Scale (TAS) and the Empathy Quotient (EQ) tools for a broader analysis, comparing performance of Llama 3.1 to GPT 3.5, GPT 4, and Gemini models. By adopting a controlled prompt-tuning method, our study explored the impact of different prompt styles, verbose and concise, and the augmenting of immediate knowledge base on the model response quality.

Keywords

Emotional Intelligence, Psychometric Assessment, Foundation Model, Prompt Engineering

Introduction

Emotional intelligence (EI) emerged as a distinct psychological research branch about three decades ago (Searle 1980). In its evolution, EI was further classified into ability, trait, and mixed principal models, each with their respective measures (Rivers et al. 2020).

The mixed EI models comprise a mixture of personality and behavioral items and measure a combination of traits, social skills, and competencies. Ability EI tests assess constructs related to theoretical understanding of emotion and are typically based on maximal performance. Whereas trait EI instruments are founded on self-report items and often utilize scales rather than yes/no answers that leads to stable psychometric properties (O'Connor et al. 2019; Bru-Luna et al. 2021). In our work, we primarily employed the widely used Trait Emotional Intelligence Questionnaire (TEIQue; Petrides 2009) tool that possesses excellent consistency and correlation with the Big Five personality model. We supplemented TEIQue by the Toronto Alexithymia Scale (TAS-20) and the Emotional Quotient Inventory (EQ-60) psychometric tests to corroborate analysis.

We explored self-assessment of the TEIQue questionnaire using the state-of-the-art Instruct version offered by the Llama 3.1 foundation model (Meta 2024). In the post-training stage, we tuned the model to follow instructions and improve answer quality by prompting the model with feedback to refine its own response. Prompt tuning (Lester et al. 2021; Li and Liang 2021) is one of the most effective solutions to reuse a frozen foundation model for a multitude of downstream tasks without retraining the model and updating all its weights. The Instruct version of the Llama 3.1 model uses a conversation structure to represent the input prompt that has to be reproduced in its entirety for commensurate performance.

Despite the demonstrated empirical efficacy of prompt tuning to adapt a pretrained transformer-based foundation model for a new task, the theoretical support of the difference between tuning parameters before the input against the tuning of model weights is limited. Recent studies (Wang et al. 2023c; Hu et al. 2024) proved statistically that prompt tuning on a simplest possible transformer architecture, comprising a single-head configuration with only a single self-attention layer are universal approximators for any sequence-to-sequence Lipschitz functions. In addition, the work validates the memorization capacity of prompt tuning and derives a lower bound on required soft-prompt tokens as exponential-in-dL and-in- $1/\epsilon$, where d is the token dimension, L the input sequence, and ϵ the approximation error. Rather than solely memorizing the last token of a pair of token sequences, Hu et al. (2024) further demonstrate a generalized memorization of prompt tuning on any general dataset.

Our analyses of knowledge transferability to the task of emotional intelligence draws to a large extent from the aforementioned essential and rigorous theoretical claims for supporting our empirical evidence. These formal theoretical derivations aid to mathematically substantiate the claim that updating only prompt embeddings L_p is sufficient for our observed performance gains. The presence of exhaustive proofs means that the validity of these theoretical claims is ultimately contingent on experimental observations and prior literature. Distinctly in our evaluation we contrast the performance impact of assessing EI by employing both succinct and verbose prompts that abide by Llama 3.1 format. Notably items in our EI questionnaires comprise a typical query length of about 12 words on average, considerably

Corresponding author:

Avi Bleiweiss

Email: avibleiweiss@bshalem.onmicrosoft,com

¹BShalem Research

¹Sunnyvale, CA, USA

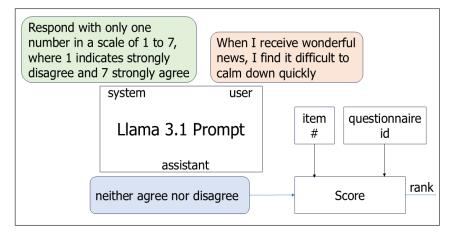


Figure 1. Overview of our proposed EI assessment framework. A questionnaire item paired with context for the query are the inputs to the Llama 3.1 user and system prompt roles, respectively. The assistant responds with a Likert seven-scale figure that follows score computation based on the questionnaire ID (TEIQue, TAS, or EQ) and an item index.

shorter than the maximal 4,096 input tokens the model permits. Thus, the token length ratio of prompt to input is relatively high in our task and increases considerably when retrieving data from an augmented knowledge base.

Background

Interpreting emotion in social context is a key element of EI for effective conversation and interaction. Despite their notable strides in a broad range of disciplines, studies on evaluating the human-like empathy traits of foundation models have been relatively scarce and confined to a single modality of textual items. Guided by the language generation process that predominately renders stateless in foundation models, EI researchers presently emphasize emotion understanding, which relies on eliciting social context from a narrative.

We briefly survey recent work on emotional intelligence assessment in large foundation models. Wang et al. (2023a) developed a novel psychometric assessment of emotion understanding following the standardized Mayer-Salovey-Caruso Emotional Intelligence Test (MSCEIT; Mayer et al. 2003). Their text-based evaluation provided consistent assessment for both humans and foundation models, however, their test comprised only forty items and renders a limited extent of traits. Paech (2024) introduced EQ-Bench that builds upon the Situational Evaluation of Complex Emotional Understanding test (SECEU; Wang et al. 2023b). Rather than a text-based query modality, the subject is presented with a dialogue and asked to rate four emotional candidates including surprised, confused, angry, and forgiving. The interpretation of the play script is critical nonetheless for obtaining plausible scores. Sabour et al. (2024) proposed a theory-based EI benchmark, EmoBench, composed of 400 human-curated items to address understanding of complex emotions. Their experiments showed that 48 screened human participants outperformed their set of current foundation models. Mozikov et al. (2024) proposed an emotion modeling framework in large foundation models by comparing LLM behavior with humans in ethical benchmarks and game-theoretical experiments. Their analysis strongly suggests that large models deviate significantly from human emotional responses.

To the extent of our knowledge, assessing EI in the instruction-tuned Llama 3.1 model using clinically-approved psychometric tests has not been explored in prior work. Our paper contributes to the increased interest in fostering emulation of human emotional traits in foundation models, and provides extensive analysis assessing their alignment with human benchmarks on three empathy understanding tests: TEIQue, TAS, and EQ. Although not generalized beyond specific EI tasks, our findings present evidence that using verbose prompts can improve performance compared to a succinct version. The effect on response quality was also analyzed by employing a retrieval-augmented generation approach to extend the immediate EI knowledge base of the foundation model. While suggesting a relative mild gain of performance, this method incurs a linear computational complexity with fetching the top-k ranked passages.

El Questionnaire

The Trait Emotional Intelligence Questionnaire (TEIQue) is founded on trait EI theory that perceives EI as a personality attribute (Petrides 2009; Andrei et al. 2016). To date, only trait EI theory offers a comprehensive scientific framework for interpreting the diverse results of independent empirical research in a way consistent with the long-standing study of individual differences in personality and emotion. TEIQue is part of a set of measures based on the trait EI model that include questionnaires for children, adolescents, and adults. TEIQue has been broadly adapted to other languages and studies relevant to its validity confirming results in line with the English version. In our evaluation, we use the most recent full-form version of TEIQue— a self-report inventory that comprises 153 items, 15 facets, four factors, and global trait EI, ¹ extending over the sampling domains of trait EI as shown in Table 4. We score the TEIQue test responses in Llama 3.1 at the questionnaire, factor, and facet levels of hierarchy.

To assess EI quality we used the seven-point Likert scale that indicate strength of agreement related to an item statement or query. Rather than a percentile measure, we report scores as real numbers $\in \{1.0, 2.0, ..., 7.0\}$ corresponding to the following answer choices: (i) strongly disagree, (ii) disagree, (iii) somewhat disagree, (iv) either agree or disagree, (v) somewhat

Bleiweiss 3

Table 1. Reviewing the Llama 3.1 prompt format in adapting the system, user, and assistant roles for assessing El item responses.

```
<|begin_of_text|>
<|start_header_id|>system<|end_header_id|> Respond with
only one number in a scale of 1 to 7, where 1 indicates
strongly disagree and 7 strongly agree<|eot_id|>
<|start_header_id|>user<|end_header_id|>I am usually
able to control other people<|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
null<|eot_id|><|start_header_id|>3<|eot_id|>
<|end_of_text|>
```

agree, (vi) agree, and (vii) strongly agree. Likert scale scores are further grouped in three measure tiers for ease of interpretation: (a) Values in the range [1.0,2.0) are considered below average, (b) Measures that span [2.0,5.0) constitute the average scale, and (c) Rates extending [5.0,7.0] are above average.

Prompt Tuning

Fine-tuning has been the prevalent approach to adapt pretrained foundation models for downstream tasks. However, fine-tuning is prohibitively computationally expensive in revising the entirety of large model weights, and requires storing a tuned copy of the model for each task. Prefix-tuning (Li and Liang 2021) and prompt-tuning (Lester et al. 2021) are forms of simplification to model tuning by freezing most of the pretrained parameters and only manage a much smaller set of task-definition parameters. All prompt tuning tasks are cast as a text generation process that provides instructions prepended to the task input text and produces the task outputs from the tuned model. Thus, given the input X, a series of n prompt tokens $P = p_{1:n}$, and the generated output Y, the foundation model maximizes the likelihood of Y, $P_{r_{\theta}}(Y|[P;X])$, while keeping the model parameters θ fixed. Prompt tuning further relaxes conditional generation by using a fixed prompt of special tokens and only the embeddings of these tokens θ_P are updated. The revised optimization function is reduced to $P_{r_{\theta;\theta_P}}(Y|[P;X])$.

Llama 3.1 is a foundation model designed to generate quality text from user inputs. Using Llama 3.1 effectively to generate guided outputs, requires a structured input format—a prompt— to interact with the model. The clarity and the context the prompt provides are essential to succeeding model responses. The Llama 3.1 prompt comprises text sequences of tuned tokens and roles that are processed by the model. In Table 1, we show a typical prompt we used for assessing EI items. The system role sets the context to the model, the user conveys an EI query or statement that constitutes a textual questionnaire item, and the assistant represents the response issued by the model: an item score of 3 that matches the somewhat-disagree answer choice.

Experiments

In our experiments, we used the instruction-tuned Llama 3.1 (Meta 2024) foundation model that renders a dense Transformer architecture of 8 Billion parameters. ² We

report our results in Likert scale with seven agreement statements ranging from strongly disagree to strongly agree and corresponding to a [1.0,7.0] scoring range. We ran inference locally and entirely on the CPU with up to four workers, while not exceeding 9.5GB of system memory. Our running time last about three minutes on average for each item in any of the test questionnaires.

Trait Hierarchy We conducted our experiments at three TEIQue hierarchical levels, including the questionnaire, factor, and facet, producing statistical distribution of EI assessments in a single global, five, and fifteen measures, respectively. We note that the association of an EI item with a facet is vital for clinical validity and is not provided explicitly in the online version of the TEIQue we obtained. Instead, with no loss of generalization to our study, we randomly generated a balanced item-facet relationship of about ten items per facet, on average. The computed mean of item responses generated in Llama 3.1 across any of the entire TEIQue questionnaire, one of the five factors, or one of the fifteen facets are compared to a corresponding human measure obtained from a large human sample (Andrei et al. 2016).

Table 2. Statistical distributions of emotional intelligence assessment for the TEIQue questionnaire.

Ouestionnaire		Human			
Questionnaire	min	max	μ	σ	μ
Global Score	1.00	7.00	4.97	1.17	4.06

Global Score In Table 2, we present a single global score of an average assessment measure at the TEIQue questionnaire level along with complementary statistical distribution of prompt assistant responses. Llama 3.1 score of 4.97 outperforms the corresponding human score of 4.06.

Table 3. Statistical distributions of emotional intelligence assessment for each TEIQue factor.

Factor		Llama 3.1			
racioi	min	max	μ	σ	μ
Emotionality	1.00	7.00	4.84	1.27	3.08
Self Control	2.00	6.00	5.04	1.20	6.02
Sociability	3.00	7.00	5.09	1.06	3.15
Trait EI	3.00	7.00	5.00	1.19	1.75
Well Being	2.00	6.00	4.97	1.15	4.83

Factor Assessment In Table 3, we show EI performance of the five TEIQue factors. The Llama 3.1 scores are consistently narrowly distributed along the upper threshold of above-average tier, and shown to exceed most of the corresponding human measures. Notably the self control factor renders a human percentile score that outperforms the foundation model score by 15 percentage points, owing to a factor most readily responsive to training. We highlight the maximal model score of 5.09 for sociability and 6.02 for human self control. Evidently Llama 3.1 surpasses human scores for all the remaining four factors by a large margin.

Facet Analysis To address the lack of an item-facet linkage in the questionnaire we obtained, we commenced the following preprocessing steps: (i) Given n items and m facets, we resized the facet vector by replication using $\lfloor n/m \rfloor$. (ii) Entries of the linked table are extended by the remainder n%m. (iii) We randomly shuffle item-facet relations in the final table. The table we generated for TEIQue, where n=153 and m=15 is fairly balanced with fourteen facets attached to ten items each and one facet with thirteen associations. Measures of the TEIQue facet assessment are shown in Table 4. Llama 3.1 surpasses human assessment measures on eleven of the fifteen facets and sustains a uniform above-average ratings across measures. Maximal facet figures of 5.64 for the model emotion regulation and 6.38 for human stress management are highlighted.

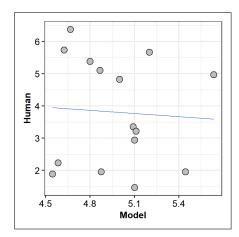


Figure 2. Scatterplot of model and human facet scores.

In Figure 2, we show a scatterplot of model and human scores for TEIQue facets. Our computed correlation coefficients are: (i) r=-0.061 for Spearman, (ii) rho=-0.093 for Pearson, and (iii) tau=-0.106 for Kendall; where confidence level is 0.95, the Spearman method uses product-moment correlation, and Pearson and Kendall apply rank correlation. The correlation coefficients we report consistently suggest a fairly weak and inverse model-human assessment relationship.

Baseline Performance Analysis To conduct a broader performance study, we administered the Toronto Alexithymia Scale (TAS-20; Bagby et al. 1994; Leising et al. 2009) and the Empathy Quotient (EQ-60; Baron-Cohen and Wheelwright 2004) questionnaires to prompt the Llama 3.1 model and score its responses in contrast to other existed foundation models. The twenty-item TAS-20 is used to measure general psychological distress comprising (i) difficulty identifying feelings (DIF), (ii) difficulty describing feelings (DDF),

and (iii) externally oriented thinking (EOT) subscales with 7, 5, and 8 items, respectively. The EQ-60 tool is a self-administered questionnaire consisting of 60 statements split into 40 empathy and 20 control items. EQ-60 is intended to instrument levels of empathy in adults with high functioning autism, considered an empathy disorder. In their study, Parker et al. (2001) used TAS-20 and EQ-60 to assess alexithymia and emotional intelligence and found that although their mutual constructs are independent, they overlap considerably and are strongly and inversely related. In the Supplemental material section, we provide details of questionnaire scoring.

In their work, Patel and Fan (2024) evaluated three current language models for their empathy levels and the capacity to identify and describe emotions by prompting them to answer questions of the TAS-20 and EQ-60 assessment tests. ³ ⁴ In Table 5, we show their global scores for GPT 3.5, GPT 4, and Gemini, contrasted with both our Llama 3.1 measures and human benchmarks scored on a 5-point Likert scale. On TAS-20, human controls administered to a clinical population sample of 1,933 adults comprising 880 men and 1,053 women with a μ_{aqe} and σ_{aqe} of 35.5 and 12.6 years, respectively (Parker et al. 2003). Noting that the lower the scores the higher the performance, the results of Llama 3.1 and Gemini suggest the presence of alexithymia at a high threshold, while GPT 3.5 renders the top score implying even a higher level of alexithymia. Conversely, GPT 4 experience no alexithymia symptoms with a score of 3.40 that is fairly aligned with human benchmarks of 3.20 on average. On EQ-60, the inverse applies, as ascending scores rather express improved performance. GPT 3.5 and GPT 4 score much lower than the neurotypical control benchmark applied to a sample of 90 adults that consisted of 65 males and 25 females with a μ_{aqe} and σ_{aqe} of 34.2 and 11.8 years, respectively (Baron-Cohen and Wheelwright 2004). Gemini, the highest empathy scoring model with 3.96 is considerably behind Llama 3.1 at 5.43, nonetheless, both surpass humans with a 2.95 score and show their ability for assessing emotional intelligence. Evidently our exploration shows that foundation models possess inherent traits that resemble human personality traits.

We also provide a lateral performance view for Llama 3.1 in assessing a global score for TEIQue (Table 2) and EQ-60 tests. The TEIQue score of 4.97 falls behind EQ-60 with 5.43, owing to an increased item complexity for TEIQue.

Discussion

Our work cast the task of assessing emotional intelligence as an optimization problem that varies the prompt textual length while seeking to improve performance. To this end, we experimented with both verbose and succinct textual versions of the system role. The scale of the prompt thus trades off simplified intelligence reasoning with challenging reading comprehension.

Prompt Token Length In the example we outline in Table 6, the verbose variant of the system role expresses all the seven Likert scales flatly, therefore founding a straightforward association with the assistant response. Whereas the much more succinct version of the prompt only provides edge information leaving the foundation model to interpret core values. The token length of the system role snippets for their verbose and succinct representations are 38 and 21,

Bleiweiss 5

Table 4. Statistical distributions of emotional intelligence assessment for each facet.

Facet	Factor		Llama 3.1			
	ractor	min	max	μ	σ	μ
Adaptability	Trait EI	3.00	6.00	4.88	1.25	1.96
Assertiveness	Sociability	3.00	6.00	5.10	1.29	2.94
Emotion Expression	Emotionality	1.00	6.00	4.55	1.75	1.89
Emotion Management	Sociability	4.00	6.00	4.87	0.92	5.11
Emotion Perception	Emotionality	4.00	7.00	5.09	1.04	3.36
Emotion Regulation	Self Control	3.00	6.00	5.64	0.92	4.97
Empathy	Emotionality	2.00	6.00	4.58	1.24	2.24
Happiness	Well Being	3.00	6.00	5.00	1.15	4.83
Impulse Control	Self Control	4.00	6.00	4.62	0.92	5.74
Self Motivation	Trait EI	4.00	7.00	5.10	1.20	1.47
Optimism	Well Being	2.00	6.00	5.11	1.36	3.22
Relationships	Emotionality	4.00	6.00	5.20	0.92	5.67
Self Esteem	Well Being	3.00	6.00	4.80	1.03	5.39
Social Awareness	Sociability	4.00	7.00	5.44	1.01	1.96
Stress Management	Self Control	2.00	6.00	4.67	1.50	<u>6.38</u>

Table 5. Baseline comparison of emotional intelligence assessment across foundation models on TAS-20 and EQ-60 questionnaires.

Study	Model	TAS-20		EQ-60	
Study	Wiodei	μ	σ	μ	σ
	GPT 3.5	5.20	0.66	2.64	0.57
Patel and Fan (2024)	GPT 4	3.40	0.44	1.94	0.38
	Gemini	4.25	0.81	3.96	0.34
Ours	Llama 3.1	4.25	0.79	5.43	0.95
Ours	Llama 3.1-rag	4.12	0.77	<u>5.47</u>	0.96
Parker et al. (2001)	Human	<u>3.20</u>	0.80	2.95	0.74

respectively, evidently well in the plausible range for prompt tuning to yield effective performance gains (Lester et al. 2021). In practice, rather than selecting prompt tokens from a vocabulary derived from a large set of frozen weights, Llama 3.1 uses for each conversation turn a fixed system role comprised of special tokens, where only the embeddings of these tokens must be updated. We note that all our thus far reported Llama 3.1 results used the compressed prompt variant.

Table 6. Verbose and compact text versions describing the system role of the prompt for a single conversational turn.

Scale	System Role
verbose	Respond with 1 if you strongly disagree, 2 if you disagree, 3 if you somewhat disagree, 4 if you either agree or disagree, 5 if you somewhat agree, 6 if you agree, and 7 if you strongly agree
compact	Respond with only one number in a scale of 1 to 7, where 1 indicates strongly disagree and 7 strongly agree

Prompt Comprehension To contrast the assessment of verbose with the succinct prompt renditions we used the TAS-20 alexithymia scale. The entirety of the twenty assistant responses to the compact system role representation consistently comprised a single numerical value between

one and seven. However, for the verbose prompt form the Llama 3.1 model emitted single numerical values for the majority of questionnaire items, but for the remaining handful the model issued phrasal utterances, such as (i) I would respond with a 2., (ii) That's a 4. You're indicating that, and (iii) You responded with a 2, disagree. Outstanding in our results is the TAS-20 total score for the verbose prompt with a 3.35 measure, compared to the compact system role variation that drew a 4.25 outcome. The verbose prompt outperforms the GPT 4 score of 3.40 (Table 5), and reached extremely close to the human assessment measure of 3.20 (Parker et al. 2003). This assessment suggests that the verbose prompt impacts performance favorably and further leading to a desireful state of no alexithymia symptoms. Llama 3.1 achieves some form of general intelligence by influencing performance on cognitive ability measures presented in our task, although the compact prompt is the less intuitive to decipher.

Retrieval-Augmented Generation (RAG) Large foundation models often produce text that lack domain-specific expertise. To mitigate this shortfall, Lewis et al. (2020) offer the inclusion of multiple top ranked passages, which provide supplementary knowledge context in LLMs and formulates a generalized fine-tuning approach for retrieval-augmented generation (RAG). RAG models integrate the parametric memory of a pretrained sequence-to-sequence architecture with a non-parametric memory— a dense vector index of knowledge-base, accessed with a pretrained neural retriever.

Thus, RAG involves incorporating a critically-reviewed external knowledge source in the LLM prompt to increase response quality. Findings of Lewis et al. (2020) study showed that RAG models generate more specific, diverse, and factual language than a state-of-the-art parametric-only sequenceto-sequence baseline. However, RAG cannot be entirely exploited during generation if the retrieved documents are not relevant. Levonian et al. (2023) evaluated RAG efficacy for math question answering (QA) using content from a highquality open-source math textbook. Their results strongly suggest that RAG improves answer quality, however, this is at the cost of excessive prompt guidance that was less favorable by the students surveyed. Similarly, we ask whether retrieval-augmented generation linked with reframing Llama 3.1 instruction prompts increases human alignment of LLM responses for EI assessment.

Given k knowledge-based passages we denote as $r_{1:k}$, $\sum_{i=1}^k |r_i|$, where $|r_i|$ is the token cardinality of a retrieved document, must be less than the maximum input token length of the foundation model— 4,096 for Llama 3.1. The passages are arranged in the order they were retrieved and combined into a single text string, further concatenated with the EI questionnaire item to form the following augmented user role: user_prompt = $(item, r_{1:k})$. This composite text is then fed into the foundation model to generate the final assistant response.

We constructed an EI knowledge-base for our RAG experiments based on the seminal article Emotional Intelligence Measures: A Systematic Review (Bru-Luna et al. 2021). 5 Fetched by the retriever and rendered as a document set, this supplementary context was further split into 36 passages— one per pdf page. Comprising a total of 17,544 tokens, the statistical distribution of words per passage has its range from 198 to 819, mean of 487, and SD of 199. On average, passage token length is of more than an order of magnitude than a questionnaire item of 12 tokens long. The passages are top-k ranked for relevance, using cosine similarity between a passage and the questionnaire items, each represented by Llama 3.1 provided embeddings. Barring truncation of the input sequence, $k \in \{5, 10, 20\}$ are typical settings for knowledge retrieval that offer evaluating performance impact with ascending k. To ensure information ranking accuracy, we ran top-k scoring on each of TAS-20 and EQ-60 questionnaires separately. The results of RAG based prompt on LLama 3.1 are shown in Table 5. RAG results using k=5 are indicating less Alexithymia on TAS-20, down from 4.26 to 4.12 and improving on Gemini scores. On EQ-60 we observed a moderate performance gain of about one percentage point. Increasing k from five to twenty and retrieving passages out of order to ensure input token count is bound under 4K proved unsubstantiated enhancement of response quality. This is probably due to the added task requiring the LLM to effectively prune passages that are less contributing, while incurring computational complexity that grows almost linearly with k.

Prompt Inconsistency In their study, (Patel and Fan 2024) used fairly non-structured prompts that vary for each LLM and each questionnaire. For example, the prompt in GPT 3.5 on TAS-20 reads: Please read each of the following statements and

carefully rate if you strongly agree, agree, neither agree nor disagree, disagree, or strongly disagree. There are no right or wrong answers or trick questions. Feeding Llama 3.1 with this prompt directly yielded different prompt output, mostly due to architectural and training differences over models. Cross-model prompt interoperability thus becomes a pivotal goal to address by transforming inputs that work well with other large language models into forms that are better optimized for Llama. Released shortly before our paper submission, Llama Prompt Ops is such a tool the research community could benefit from and sustain consistent prompt output by abstracting the underlying LLM. ⁶

Conclusion

The key contributions of our paper are deeply intertwined with broader scientific literature on emotional intelligence in computational models and the emerging field of prompt engineering. Building on prior research that has applied psychometric tools to evaluate human-like cognition in AI, such as studies using the TEIQue, TAS-20, and EQ-60 to quantify emotional awareness, this work advances the discussion by demonstrating how prompt tuning can enhance the fidelity of these assessments in large language models. It aligns with findings in the literature that emphasize the importance of context-rich inputs in generating higher quality, more consistent outputs, while also resonating with recent developments in parameter-efficient fine-tuning techniques. Moreover, by comparing Llama 3.1 with other notable model like GPT 3.5, GPT 4, and Gemini, the paper extends previous comparative analyses in the field, highlighting both convergent and divergent performance trends that enrich the dialogue on model adaptation and evaluation in AI research.

We experimented with system-level prompts to provide Llama 3.1 the most simple and succinct context to assess EI questionnaire items efficiently. We made the prompt transferable across three EI assessment tools, and showed Llama 3.1 outscore human benchmarks on most accounts. To expand on our work, an apparent venue to pursue is in interpreting complex emotion from multimodal input cues, such as text combined with facial expressions or the tone of voice.

Limitations

To ensure a stable and robust analysis, we avoided assessing neither TAS-20 subscales and nor EQ-60 factors due to their small sample size, and only report questionnaire global scores. Randomly creating item-facet relationships in TEIQue could possibly introduce facet score mismatches between foundation models and human benchmarks. To mitigate this shortcoming of reproducibility, we plan to make the code for generating these relationships publicly accessible. Evaluating TEIQue in closed-source foundation models is outside the scope of this study, and contrasting EI assessment of open-source models, Llama 3.1 and Gemini, is only conducted on TAS-20 and EQ-60 questionnaires.

Bleiweiss 7

Impact Statement

We honor and support the Sage Code of Ethics. The goal of our study is to advance the field of machine learning from the perspective of empathy assessment and its alignment with humans. Our questionnaire data was scraped from the internet; forms, versions, and translations for each of TEIQue, TAS-20, and EQ-60 are free of charge for academic research. Our study refrains from crowdsourcing or any form of direct research with human subjects. The only exception is in comparing emotional intelligence measures between the Llama 3.1 foundation model and an average profile assessment of a large sample of humans, conducted by an external study (Parker et al. 2003). We suppose that by the presence of user consent this on its own bears no plausible negative societal impact.

Supplemental material

Questionnaire Scoring We follow with a summary of the rules for scoring the TAS-20 and EQ-60 dataset.

TAS-20 There are five items that are negatively keyed: (4,5,10,18,19). They are reverse-coded, meaning that a higher score indicates a lower level of alexithymia.

EQ-60 Two points are scored for each of the following items if the answer is definitely agree, or one point if the answer is slightly agree: (1, 6, 19, 22, 25, 26, 35, 36, 37, 38, 41, 42, 43, 44, 52, 54, 55, 57, 58, 59, 60).

Conversely, two points are scored for each of the following items if the answer is definitely disagree, or one point if answered slightly disagree: (4, 8, 10, 11, 12, 14, 15,18, 21, 27, 28, 29, 32, 34, 39, 46, 48, 49, 50).

All other questions: (2, 3, 5, 7, 9, 13, 16, 17, 20, 23, 24, 30, 31, 33, 40, 45, 47, 51, 53, 56) are not scored.

Acknowledgements

We would like to thank the anonymous reviewers for their insightful suggestions and feedback.

Notes

- 1. https://www.eiconsortium.org/measures/teique.html
- https://huggingface.co/meta-llama/Meta-Llama-3.
 1-8B-Instruct
- 3. https://contextualscience.org/TAS_Measure
- 4. https://psychology-tools.com/test/empathy-quotient
- https://pmc.ncbi.nlm.nih.gov/articles/PMC8701889/pdf/ healthcare-09-01696.pdf
- 6. https://github.com/meta-llama/llama-prompt-ops

References

- Andrei F, Siegling AB, Aloe AM, Baldaro B and Petrides KV (2016) The incremental validity of the trait emotional intelligence questionnaire (TEIQue): A systematic review and meta-analysis. *Journal of Personality Assessment* 98(3): 261–276. DOI: 10.1080/00223891.2015.1084630.
- Bagby R, Parker JD and Taylor GJ (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. DOI:https://doi.org/10.1016/0022-3999(94)90005-1.

Baron-Cohen S and 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. DOI: https://doi.org/10.1023/b:jadd.0000022607.19833.00.

- Bru-Luna LM, Martí-Vilar M, Merino-Soto C and Cervera-Santiago JL (2021) Emotional intelligence measures: A systematic review. *Healthcare* 9(12): 1696. DOI:https://doi.org/10.3390/ healthcare9121696.
- Hu JYC, Wang WP, Gilani A, Li C, Song Z and Liu H (2024) Fundamental limits of prompt tuning transformers: Universality, capacity and efficiency. URL https://arxiv.org/abs/2411.16525.
- Leising D, Grande T and Faber R (2009) The toronto alexithymia scale (tas-20): A measure of general psychological distress. *Journal of Research in Personality* 43(4): 707–710. DOI: https://doi.org/10.1016/j.jrp.2009.03.009.
- Lester B, Al-Rfou R and Constant N (2021) The power of scale for parameter-efficient prompt tuning. In: Moens MF, Huang X, Specia L and Yih SWt (eds.) *Empirical Methods in Natural Language Processing (EMNLP)*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, pp. 3045–3059. DOI:10.18653/ v1/2021.emnlp-main.243.
- Levonian Z, Li C, Zhu W, Gade A, Henkel O, Postle ME and Xing W (2023) Retrieval-augmented generation to improve math question-answering: Trade-offs between groundedness and human preference. In: NeurIPS'23 Workshop on Generative AI for Education (GAIED). DOI:10.1109/ICKG63256.2024. 00023.
- Lewis P, Perez E, Piktus A, Petroni F, Karpukhin V, Goyal N, Küttler H, Lewis M, Yih Wt, Rocktäschel T, Riedel S and Kiela D (2020) Retrieval-augmented generation for knowledge-intensive nlp tasks. In: Larochelle H, Ranzato M, Hadsell R, Balcan M and Lin H (eds.) *Advances in Neural Information Processing Systems*, volume 33. Curran Associates, Inc., pp. 9459–9474. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf.
- Li XL and Liang P (2021) Prefix-tuning: Optimizing continuous prompts for generation. In: Zong C, Xia F, Li W and Navigli R (eds.) *Association for Computational Linguistics and (ACL)*. Online: Association for Computational Linguistics, pp. 4582–4597. DOI:10.18653/v1/2021.acl-long.353.
- Mayer JD, Salovey P, Caruso DR and Sitarenios G (2003) Measuring emotional intelligence with the MSCEIT v2.0. *Emotion* (*Washington, D.C.*) 3(1): 97—105. DOI:10.1037/1528-3542.3. 1.97.
- Meta (2024) The llama 3 herd of models. Technical report, Meta. https://arxiv.org/abs/2407.21783.
- Mozikov M, Severin N, Bodishtianu V, Glushanina M, Nasonov I, Orekhov D, Pekhotin V, Makovetskiy I, Baklashkin M, Lavrentyev V, Tsvigun A, Turdakov D, Shavrina T, Savchenko A and Makarov I (2024) EAI: Emotional decision-making of LLMs in strategic games and ethical dilemmas. In: *Neural Information Processing Systems*. URL https://openreview.net/forum?id=8aAaYEwNR4.
- O'Connor PJ, Hill A, Kaya M and Martin B (2019) The measurement of emotional intelligence: A critical review of the literature and recommendations for researchers and practitioners. *Frontiers in Psychology* 10(1116). DOI:10.3389/fpsyg.2019.01116.

Paech SJ (2024) Eq-bench: An emotional intelligence benchmark for large language models. URL https://arxiv.org/abs/2312.06281.

- Parker JD, Taylor GJ and Bagby R (2001) The relationship between emotional intelligence and alexithymia. *Personality and Individual Differences* 30(1): 107–115. DOI:https://doi.org/10.1016/S0191-8869(00)00014-3.
- Parker JD, Taylor GJ and Bagby R (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:https://doi.org/10.1016/S0022-3999(02) 00578-0.
- Patel SC and Fan J (2024) Identification and description of emotions by current large language models. *bioRxiv* DOI:10.1101/2023. 07.17.549421. URL https://www.biorxiv.org/content/early/2024/07/13/2023.07.17.549421.
- Petrides KV (2009) *The Trait Emotional Intelligence Questionnaire* (*TEIQue*). London Psychometric Laboratory, London, UK.
- Rivers SE, Handley-Miner IJ, Mayer JD and Caruso DR (2020) Emotional intelligence. In: Sternberg RJ (ed.) *The Cambridge handbook of intelligence*, 2nd edition. Online: Cambridge University Press, pp. 709–735.
- Sabour S, Liu S, Zhang Z, Liu J, Zhou J, Sunaryo A, Lee T, Mihalcea R and Huang M (2024) EmoBench: Evaluating the emotional intelligence of large language models. In: Ku LW, Martins A and Srikumar V (eds.) Annual Meeting of the Association for Computational Linguistics (ACL). Bangkok, Thailand: Association for Computational Linguistics, pp. 5986– 6004. URL https://aclanthology.org/2024.acl-long.326.
- Searle JR (1980) Minds, brains, and programs. *Behavioral and Brain Sciences* 3(3): 417–424. DOI:10.1017/S0140525X00005756.
- Wang X, Li X, Yin Z, Wu Y and Liu J (2023a) Emotional intelligence of large language models. *Journal of Pacific Rim Psychology* 17: 18344909231213958. DOI:10.1177/18344909231213958.
- Wang X, Li X, Yin Z, Wu Y and Liu J (2023b) Emotional intelligence of large language models. *Journal of Pacific Rim Psychology* 17. DOI:10.1177/18344909231213958.
- Wang Y, Chauhan J, Wang W and Hsieh CJ (2023c) Universality and limitations of prompt tuning. In: Oh A, Naumann T, Globerson A, Saenko K, Hardt M and Levine S (eds.) Advances in Neural Information Processing Systems, volume 36. Curran Associates, Inc., pp. 75623–75643. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/ eef6aecfe050b556c6a48d9c16b15558-Paper-Conference.pdf.