Large Language Models Can Infer Psychological Dispositions of Social Media Users

Heinrich Peters

Columbia University New York hp2500@columbia.edu Sandra Matz

Columbia University New York sm4409@columbia.edu

ABSTRACT

Large Language Models (LLMs) demonstrate increasingly human-like abilities across a wide variety of tasks. In this paper, we investigate whether LLMs like ChatGPT can accurately infer the psychological dispositions of social media users and whether their ability to do so varies across socio-demographic groups. Specifically, we test whether GPT-3.5 and GPT-4 can derive the Big Five personality traits from users' Facebook status updates in a zero-shot learning scenario. Our results show an average correlation of r = .29 (range = [.22, .33]) between LLM-inferred and self-reported trait scores – a level of accuracy that is similar to that of supervised machine learning models specifically trained to infer personality. Our findings also highlight heterogeneity in the accuracy of personality inferences across different age groups and gender categories: predictions were found to be more accurate for women and younger individuals on several traits, suggesting a potential bias stemming from the underlying training data or differences in online self-expression. The ability of LLMs to infer psychological dispositions from usergenerated text has the potential to democratize access to cheap and scalable psychometric assessments for both researchers and practitioners. On the one hand, this democratization might facilitate large-scale research of high ecological validity and spark innovation in personalized services. On the other hand, it also raises ethical concerns regarding user privacy and self-determination, highlighting the need for stringent ethical frameworks and regulation.

Keywords Large language models · ChatGPT · GPT-4 · Personality · Big Five

1 Introduction

Large language models (LLMs) and other transformer-based neural networks have revolutionized text analysis in research and practice. Models such as OpenAI's GPT-4 [1] or Anthropic's Claude [2], for example, have shown a remarkable ability to represent, comprehend, and generate human-like text. Compared to prior NLP approaches, one of the most striking advances of LLMs is their ability to generalize their "knowledge" to novel scenarios, contexts, and tasks [3, 4].

While LLMs were not explicitly designed to capture or mimic elements of human cognition and psychology, recent research suggests that – given their training on extensive corpora of human-generated language – they might have spontaneously developed the capacity to do so. For example, LLMs display properties that are similar to the cognitive abilities and processes observed in humans, including theory of mind (i.e., the ability to understand the mental states of other agents [5]), cognitive biases in decision-making [6] and semantic priming [7]. Similarly, LLMs are able to effectively generate persuasive messages tailored to specific psychological dispositions (e.g., personality traits, moral values [8]).

Here, we examine whether LLMs possess another quality that is fundamentally human: The ability to "read" people and form first impressions about their psychological dispositions in the absence of direct or prior

interaction. As research under the umbrella of zero-acquaintance studies shows, people can be remarkably accurate at judging the psychological traits of strangers simply by observing traces of their behavior under certain conditions [9]. While such judgments can be influenced by stereotypes and their accuracy can vary based on the traits being assessed and the context in which judgments are made [10], past work indicates that people are able to predict a stranger's personality traits by observing their offices or bedrooms [11], examining their music preferences [12], or scrolling through their social media profiles [13].

Existing research in computational social science shows that supervised machine learning models are able to make similar predictions. That is, given a large enough dataset including both self-reported personality traits and people's digital footprints - such as Facebook Likes, music playlists, or browsing histories – machine learning models are able to statistically relate both inputs in a way that allows them to predict personality traits after observing a person's digital footprints [14, 15]. This is also true for various forms of text data, including social media posts [16, 17], personal blogs [18], or short text responses collected in the context of job applications [19].

In this paper, we test whether LLMs have the ability to make similar psychological inferences without having been explicitly trained to do so (known as zero-shot learning [3]). Specifically, we use Open AI's ChatGPT (GPT-3.5 and GPT-4 [1]) to explore whether LLMs can accurately infer the Big Five personality traits Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism [20] of social media users from the content of their Facebook status updates in a zero-shot scenario. In addition, we test for biases in ChatGPT's judgments that might arise from its foundation in equally biased human-generated data. Building on previous work highlighting inherent stereotypes in pre-trained NLP models [21, 22], we explore the extent to which the personality inferences made by ChatGPT are indicative of gender and age-related biases (e.g., potential biases in how the personality of men and women or older and younger people is judged).

Understanding the capabilities and limitations of LLMs with regard to inferring highly intimate psychological traits from digital footprints is critical, given their rapid adoption in both research and practice. On the one hand, easy access to the psychological profiles of individuals creates unprecedented opportunities to study individual differences at scale and customize products, services, or behavioral interventions to individuals' unique dispositions. On the other hand, however, automated psychological inferences pose considerable ethical and legal challenges with regard to individuals' privacy and self-determination [23]. This problem is exacerbated by the fact that LLMs are also able to automatically craft persuasive messages based on users' personality profiles [8]. The combination of fully automated psychological assessments and personalized interactions opens the door for manipulation and misuse at scale and with little to no human oversight. We discuss our findings in light of both the opportunities for scientists and practitioners and the challenges that will require new forms of AI governance and regulation [24–26].

2 Method

2.1 Data and Sampling

Our analyses are based on text data obtained from MyPersonality [27], a Facebook application that allowed users to take real psychometric tests - including a validated measure of the Big Five personality traits (IPIP [28]) - and receive immediate feedback on their responses. Users also had the opportunity to donate their Facebook profile information - including their public profiles, Facebook Likes, and status updates – to research. For the purpose of this study, we randomly subsampled 1,000 adult users $(24.2 \pm 8.8 \text{ years old}, 63.1\% \text{ female})$ who completed the full 100-item IPIP personality questionnaire and had at least 200 Facebook status updates (if they had more, we used the most recent 200). The study received IRB approval from Columbia University's ethics review board (Protocol #AAAU8559).

2.2 Measures

MyPersonality measured users' personality traits using the International Personality Item Pool (IPIP [28]), a widely established self-report questionnaire that captures the Big Five personality traits of Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism [20]. We only included users who had completed the full questionnaire with all 100 items.

To obtain inferred personality traits from ChatGPT, we used the last 200 Facebook status updates generated by each user without additional preprocessing. The average length of status updates in our sample was 17.10 words (SD=15.03). Status updates were scored using the ChatGPT API with GPT-3.5 (version gpt-3.5-turbo-0301) and GPT-4 (version gpt-4-0314) [1] as underlying models. For this purpose, the status updates were

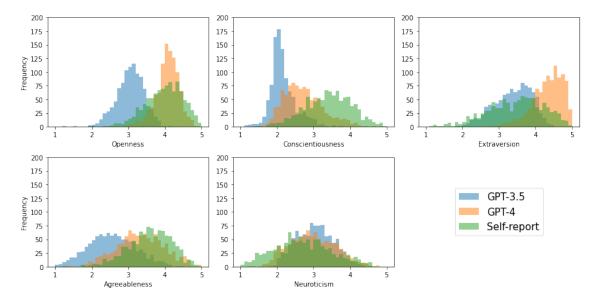


Figure 1: Distributions of self-reported and inferred personality scores for GPT-3.5 and GPT-4. Histograms show absolute frequencies for an overall sample size of n=1000. GPT-3.5 underestimates Openness. Both models underestimate Conscientiousness and Agreeableness but overestimate Neuroticism. For Extraversion, the two models diverge with GPT-3.5 underestimating and GPT-4 overestimating the true scores. Overall, GPT-4 inferred scores were more aligned with self-reported scores, indicating a potential improvement over GPT-3.5.

first concatenated into chunks and then fed into the GPT model, using a set of simple prompts to guide the behavior of the model. The system prompt was the default for GPT-3.5 and GPT-4, respectively: "You are a helpful assistant". Additionally, we prompted the model to infer Big Five traits using the inference prompt: "Rate the text on the Big Five personality dimensions. Pay attention to how people's personalities might be reflected in the content they post online. Provide your response on a scale from 1 to 5 for the traits Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism. Provide only the numbers." We then used a simple text-parsing script to transform the responses into numerical scores. In order to avoid exceeding the GPT token limit, status update histories were processed in chunks of 20 messages, and the inferred personality scores were then averaged to derive overall scores.

To boost the reliability of the inferred personality estimates, we queried ChatGPT three times for each inference. Agreement across ratings across rating rounds was high for all traits (Openness: $r_{GPT3.5}$ = .88, r_{GPT4} = 0.73; Conscientiousness: $r_{GPT3.5}$ = .88, r_{GPT4} = 0.91; Extraversion: $r_{GPT3.5}$ = .92, r_{GPT4} = 0.87; Agreeableness: $r_{GPT3.5}$ = .96, r_{GPT4} = 0.94; Neuroticism: $r_{GPT3.5}$ = .91, r_{GPT4} = 0.93), and all p-values were smaller than .001 with Bonferroni correction for multiple comparisons. Given the high level of agreement, we computed aggregate inferred scores by averaging scores across the three rounds of rating. We used the aggregate scores for all further analyses.

3 Results

3.1 Can LLMs Infer Personality Traits From Social Media Posts?

In order to assess the capacity of LLMs to infer psychological traits from social media data, we compared the inferred Big Five personality scores with self-reported scores. A comparison of the distributions suggests that both versions of ChatGPT tended to underestimate Conscientiousness and Agreeableness while overestimating Neuroticism. For Openness and Extraversion, the deviations were inconsistent across ChatGPT versions: While GPT-3.5 tended to underestimate Openness and Extraversion, GPT-4 tended to overestimate Extraversion. Overall, the distributions of inferred scores were more closely aligned with self-reported scores for GPT-4 compared to GPT-3.5, suggesting a potential improvement across versions (see Figure 1). Detailed descriptive statistics can be found in S1.

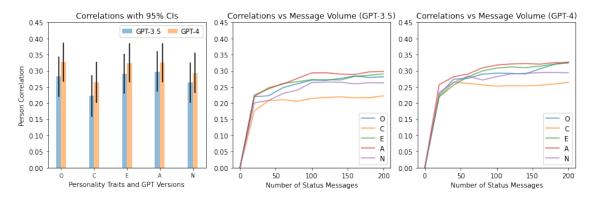


Figure 2: Pearson's correlation coefficients between inferred and self-reported scores with 95% confidence intervals (left), and Pearson's correlation coefficients for GPT-3.5 (mid) and GPT-4 as a function of message volume (right). O: Openness; C: Conscientiousness; E: Extraversion; A: Agreeableness; N: Neuroticism. Inferences for Openness, Extraversion, and Agreeableness were more accurate than those for Conscientiousness and Neuroticism, but the differences remained non-significant. Higher message volume was associated with higher levels of predictive accuracy, but a substantial share of variance was captured in as little as 20 status messages.

Importantly, the mere comparison of distributions does not provide insights into the strength and directionality of the relationships between inferred and self-reported scores. For this purpose, we conducted correlation analyses. The average Pearson correlation coefficient of inferred and self-reported scores across all personality traits was $r_{GPT3.5}$ =0.27 and r_{GPT4} =0.31. The correlations were highest for the traits of Openness ($r_{GPT3.5}$ =.28; r_{GPT4} =.33), Extraversion ($r_{GPT3.5}$ =.29; r_{GPT4} =.32) and Agreeableness ($r_{GPT3.5}$ =.30, r_{GPT4} =.32), and were slightly lower for Conscientiousness ($r_{GPT3.5}$ =.22; r_{GPT4} =.26) and Neuroticism ($r_{GPT3.5}$ =.26; r_{GPT4} =.29). All correlation coefficients were significantly different from 0 at p < .001 with Bonferroni correction for multiple comparisons. Similar to the comparison of distributions, GPT-4 showed higher levels of accuracy across all five personality traits, although none of the individual comparisons reached statistical significance (see Figure 2). Detailed results, including confidence intervals and significance levels, can be found in S2.

In addition to exploring the capacity of ChatGPT to infer personality traits from social media user data, we also tested the extent to which this capacity is sensitive to changes in the amount of data that was available for inference. Specifically, we computed correlations between self-reported and inferred personality scores based on different numbers of status messages. Specifically, we computed correlations obtained from inferences for a single chunk of status messages (20 status messages) all the way up to ten chunks (200 status messages). As expected, having access to more status messages resulted in more accurate inferences. Notably, however, most correlations are close to their maximum level after observing far less than the ultimate number of 200 status messages. In addition, the inference of certain traits seems to be particularly susceptible to the volume of input data. For example, the models' accuracy kept increasing with higher levels in input volume for Openness, Extraversion, Agreeableness, and Neuroticism, while the benefits of additional status messages leveled off earlier for Conscientiousness. See Figure 2 for a graphical representation and S3 for detailed statistics.

3.2 Does the Quality of LLM Inferences Vary Across Demographic Groups?

In order to uncover potential gender and age-related biases, we analyzed group differences in inferred Big Five scores, as well as their residuals with respect to self-reported scores. Notably, such gender and age differences might not only emerge in inferred personality scores but are also known to exist in self-reports [29, 30]. Consequently, we test for both overall group differences and differences in the residuals between the self-reported and inferred personality scores of each individual.

3.2.1 Gender Differences

We first explored the extent to which any observed group differences in inferred personality traits across men and women aligned with those observed in self-reports. As Figure 3 shows, women tend to score

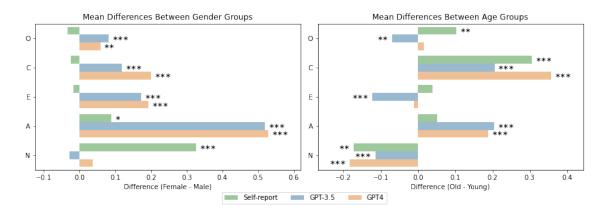


Figure 3: Mean differences in personality scores between gender groups (left) and age groups (right) for self-reported scores as well as inferences by GPT-3.5 and GPT-4. Positive values indicate higher scores for female users compared to male users and older users compared to younger users. O: Openness; C: Conscientiousness; E: Extraversion; A: Agreeableness; N: Neuroticism. ***p<.001; **p<.01; *p<.05. The results show significant gender and age differences across all personality traits.

significantly higher in Agreeableness (t=2.31; p=.021) and Neuroticism (t.=6.53; p<.001) when these traits are measured using questionnaires. In contrast, women scored significantly higher in Openness (t_{GPT3.5}= 3.42, p_{GPT3.5}< .001; t_{GPT4}= 2.72, p_{GPT4}= .007), Conscientiousness (t_{GPT3.5}= 5.28, p_{GPT3.5}< .001; t_{GPT4}= 5.73, p_{GPT4}< .001), Extraversion (t_{GPT3.5}= 5.21, p_{GPT3.5}< .001; t_{GPT4}= 7.25, p_{GPT4}< .001), and Agreeableness (t_{GPT3.5}= 13.53, p_{GPT3.5}< .001; t_{GPT4}= 13.63, p_{GPT4}< .001) when these traits were inferred by ChatGPT models, with no significant differences found for Neuroticism. This finding offers initial evidence for potential gender biases in the personality inferences made by LLMs (see Figure 3).

To further explore these potential biases, we analyzed the residuals between inferred scores and self-reported scores as an indication of how well GPT is able to represent the personality traits of male and female users. The findings suggest that GPT's personality inferences are less accurate for men than women. First, we observed larger absolute residuals for male users in Conscientiousness ($t_{GPT3.5}$ = -3.53, $p_{GPT3.5}$ < .001; t_{GPT4} = -4.48, p_{GPT4} < .001), Agreeableness ($t_{GPT3.5}$ = -9.22, $t_{GPT3.5}$ < .001; t_{GPT4} = -5.22, t_{GPT4} < .001), and Neuroticism ($t_{GPT3.5}$ = -4.55, $t_{GPT3.5}$ < .001; t_{GPT4} = -2.39, t_{GPT4} = .017) across both GPT models, indicating lower accuracy on these traits for men. Additionally, we found larger residuals for male users for GPT-3.5 in Openness ($t_{GPT3.5}$ = -3.84, $t_{GPT3.5}$ < .001; t_{GPT4} = -0.92, t_{GPT4} = .357) and larger residuals for female users in Extraversion for GPT-4 ($t_{GPT3.5}$ = -1.36, $t_{GPT3.5}$ < .173; t_{GPT4} = 3.12, t_{GPT4} = .002). For a visual representation, please refer to Figure 4). Detailed statistics can be found in S4.

Taken together, the findings suggest that GPT's personality inferences are less accurate for men than women. Notably, however, these biases seem to be limited to the absolute measures of accuracy and do not necessarily translate to GPT's ability to make inferences about men's relative personality levels. That is, when computing Pearson correlations within gender groups, we did not observe any significant difference in the magnitude of these correlations. Similarly, controlling for gender in the overall correlations between self-reported and inferred personality scores by z-standardizing inferred scores within each gender group did not yield correlations significantly different from those obtained before.

3.2.2 Age Differences

As for gender, we first explored the extent to which any observed group differences in inferred personality traits across younger and older adults (classified using a median split) were aligned with those observed in self-reports. As Figure 3 shows, older users displayed significantly higher self-reported scores in Openness (t=2.96; p=.003) and Conscientiousness (t=7.27; p<.001) and significantly lower self-reported scores in Neuroticism (t=-3.28; p=.001) compared to younger users. Partially mimicking these differences in self-reported personality traits, inferred scores were significantly higher in Conscientiousness ($t_{GPT3.5}$ = 9.23, $p_{GPT3.5}$ < .001; t_{GPT4} = 10.41, p_{GPT4} < .001) and Agreeableness ($t_{GPT3.5}$ = 4.87, $p_{GPT3.5}$ < .001; t_{GPT4} = 4.39, p_{GPT4} < .001), and lower in Neuroticism ($t_{GPT3.5}$ = -3.43, $p_{GPT3.5}$ < .001; t_{GPT4} = -4.37, p_{GPT4} < .001) for older compared to younger users. For Openness ($t_{GPT3.5}$ = -2.86, $t_{GPT3.5}$ = .004; t_{GPT4} = -0.72,

 p_{GPT4} = .472) and Extraversion ($t_{GPT3.5}$ = -3.55, p<.001; t_{GPT4} = 0.36, p_{GPT4} = .717), older individuals scored significantly lower on inferred scores for GPT-3.5 but not GPT-4 (see Figure 3).

As before, we further explore these differences by analyzing age differences in the residuals between self-reported and inferred scores. Unlike in the analyses of gender, we found substantial inconsistency in the group differences between GPT-3.5 and GPT-4. While the inferences made by GPT-3.5 showed significantly larger absolute residuals for older users in Openness ($t_{GPT3.5}$ = 4.78, $p_{GPT3.5}$ < .001), Conscientiousness ($t_{GPT3.5}$ = 2.64, $p_{GPT3.5}$ = .008) and smaller residuals for Agreeableness ($t_{GPT3.5}$ = -2.64, $t_{GPT3.5}$ = .009), no differences in absolute residuals were found for GPT-4. For a visual representation, please refer to Figure 4) Detailed statistics can be found in S5.

Taken together, the findings suggest that ChatGPT's personality inferences might be less accurate for older adults. However, as before, these biases did not translate to ChatGPT's ability to make inferences about people's relative personality levels. We did not find significant differences between within-group correlation coefficients, and z-standardizing personality scores within age groups did not yield correlation coefficients significantly different from those reported before.

3.3 Agreement With Third-Person Observer Ratings

We conducted a preliminary analysis examining the correlations between self-reported personality scores and third-person observer ratings, as well as between LLM-inferred scores and third-person observer ratings. This allowed us to 1) compare the quality of LLM inferences against a strong human benchmark (i.e. ratings from people who have access to more identity cues than just social media profiles), and 2) examine the level of agreement between LLM inferences and human judgments. Third-person ratings were collected by letting users' Facebook friends complete a 10-item version of the IPIP personality questionnaire [28, 31] about them. The analysis includes a subset of 68 individuals for whom third-person ratings were available.

The results show that correlations between self-reported scores and observer ratings ranged from r=.198 to r=.378 (mean=.304), while the correlations between LLM-inferred scores and observer ratings ranged from r=.057 to r=.457 (mean=.269) for GPT-3.5 and r=.152 to r=.400 (mean=.276) for GPT-4 (see S6 for detailed results). Overall, the correlation coefficients were largely in the same range as those between self-reported and LLM-inferred scores. The analyses thus suggest that the accuracy of LLM inferences is on par with that of human observers. They also suggest that the LLM is using similar cues to human judges. This is true even though these cues may not always be valid predictors of people's self-perceptions. For instance, in the case of Conscientiousness, the agreement between LLM inferences and observer ratings was higher than the agreement of either the LLM or human observers with participants' self-reports.

4 Discussion

4.1 Interpretation of Results

Our findings suggest that LLMs, such as ChatGPT, can infer psychological dispositions from people's social media posts without having been explicitly trained to do so. They also offer preliminary evidence that LLMs might generate more accurate inferences for women and younger individuals (compared to men and older adults). Notably, the overall accuracy of the observed inferences (Pearson correlations between self-reported and inferred personality traits ranging between r = .22 and .33, average = .29) is slightly lower than that accomplished by supervised models which have been trained or fine-tuned specifically for this purpose and with the same textual data source as used in testing (e.g., Park et al. [17], who reported correlations between r = .26 and r = .41, average r = .37). Yet, the ability of LLMs to produce inferences of reasonably high accuracy in zero-shot learning scenarios has both important theoretical and practical implications.

Our study contributes to a growing body of research comparing the abilities of LLMs to those observed in humans [5, 7, 8]. As our findings suggest, LLMs might have the human-like ability to "profile" people based on their behavioral traces, without ever having had direct interactions with them. Although most social media posts do not contain explicit references to a person's character, ChatGPT – just like human judges [13, 32] or supervised models [31] – is able to translate people's accounts of their daily activities and preferences into a holistic picture of their psychological dispositions. Our results are aligned with previous work suggesting that Openness and Extraversion are more easily inferred than other traits [13, 31]. At the same time, LLM inferences were more congruent with observer ratings than self-reports in the case of Conscientiousness, indicating that LLMs may also replicate biases in human judgment for certain traits.

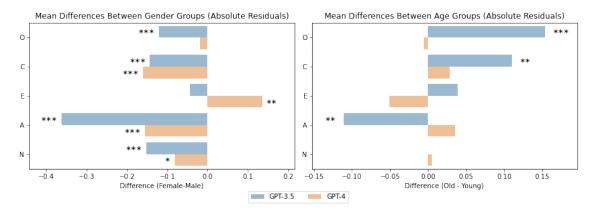


Figure 4: Mean differences in absolute residuals between gender groups (left) and age groups (right) for inferences by GPT-3.5 and GPT-4. Positive values indicate higher residuals for female users compared to male users and older users compared to younger users. O: Openness; C: Conscientiousness; E: Extraversion; A: Agreeableness; N: Neuroticism. ***p<.001; **p<.05. The results indicate lower residuals for female users in all personality traits except Extraversion. Age-related biases were observed for Openness, Conscientiousness, and Agreeableness in inferences by GPT-3.5 but not GPT-4.

Notably, the specific pathways by which LLMs such as ChatGPT arrive at their judgments and the reasons for why certain biases are introduced into the predictions (e.g., systematic gender and age differences) remain unknown. That is, we cannot speak to the question of whether LLMs use the same behavioral cues as humans or supervised machine learning models when translating behavioral residues into psychological profiles or offer an in-depth explanation for the observed differences in accuracy across age and gender categories. For example, the fact that ChatGPT shows systematic biases in its estimation of certain personality traits and is more accurate for women and younger adults could either be indicative of a bias introduced in the training of the models and/or the corpora of text data the models have been trained on, or reflective of differences in people's general self-expression on social media.

Specifically, past work indicates that LLMs are susceptible to stereotyping and bias with regard to demographic and geographic groups [21, 33–37], likely reflecting groups' representation in the underlying training data. At the same time, past work has shown differences in social media use and online self-expression across demographic groups, including age and gender [38–42]. While the past literature does not directly speak to differences in personality expression, the observed pattern of results would indicate that women and younger individuals tend to reveal more accurate information about their personalities online.

4.2 Limitations and Future Research

Our study has several limitations that should be addressed by future research. Firstly, as mentioned above, the black box character of LLMs prevents us from examining the precise mechanisms by which personality inferences are derived. As a first step in this direction, future research should analyze cue utilization by investigating which language features are highly correlated with inferred trait scores. Similarly, it would be useful to examine which language features are predictive of inference errors in order to better understand the origins of the observed gender and age biases.

Second, the text data used in our analysis was obtained from the MyPersonality Facebook application [27], which was active between 2007 and 2012. Linguistic conventions from this period might differ from contemporary online language, potentially limiting the zero-shot performance of LLMs, which have been trained on newer data. As a result, we would expect the personality inferences of LLMs to be even more accurate when applied to more contemporary data.

Third, our data was sourced from Facebook users who interacted with the MyPersonality application. As such, our sample might not be representative of the broader population of social media users (or people more generally), which could limit the external validity of our findings. For example, the general underestimation of personality traits such as Openness might be due to the fact that myPersonality users were particularly curious and open-minded.

Fourth, while our study probed how sensitive the accuracy of LLM-based inferences is to the volume of text input, we limited our data to the 200 most recent status updates. In practice, predictive performance might vary for users with fewer or more status updates. Relatedly, due to the inherent token limit in models like ChatGPT, all input data was processed in chunks. It is possible that the accuracy of future models with the ability to process larger amounts of input data at once might be higher.

Fifth, our study did not encompass the dynamics of live interactions between LLMs and users. Real-time interactions might yield different insights and highlight additional complexities not captured in our static textual data set [43]. Relatedly, while our research underscores the potential for LLMs in personalizing interactions and enhancing social computing, it does not examine the specifics of how these personalizations can be effectively implemented.

Sixth, the current research demonstrates the potential of out-of-the-box LLMs for inferring psychological variables using simple techniques such as zero-shot learning and commercially available models. It is likely that the predictive performance of LLMs could be improved through more sophisticated prompting strategies, such as chain-of-thought prompting [44] and a combination of in-context learning and supervised fine-tuning [45]. While we purposefully focused on zero-shot learning in order to establish a lower bound of predictive accuracy and investigate LLMs' inherent ability to make such predictions, future research could focus on identifying levers that maximize predictive accuracy. Aside from more sophisticated prompting paradigms, this could include giving LLMs access to users' demographic information which is typically available to human perception and could moderate the interpretation of personality-related signals. For example, the content of status messages may be interpreted differently depending on whether the user is an 18-year-old man or a 55-year-old woman. Being able to interpret message content in the context of sender identity could lead to improved inferences, but could also amplify implicit biases that are known to persist in language models [21, 33, 35].

Finally, while we make an effort to discuss the societal implications of our findings (see below), detailed recommendations regarding privacy concerns and the potential for misuse should be addressed in future research.

4.3 Implications

Our findings also have important practical implications for the application of automated psychological profiling in research and industry. Specifically, the ability of LLMs to infer psychological traits from social media data could foreshadow a remarkable shift in the accessibility - and therefore potential use – of scalable psychometric assessments. For decades, the assessment of psychological traits relied on the use of self-report questionnaires, which are known to be prone to self-report biases and difficult to scale due to their costly and time-consuming nature [46]. With the introduction of automated psychological assessments driven by supervised machine learning models [14, 19], scientists and practitioners were afforded an alternative approach that promised to expand the study and application of individual differences to research questions and domains that were previously impractical if not impossible (e.g., the use of personality traits in targeted advertising [47]; or the investigation of individual differences in large scale, ecologically valid observational studies [48]). However, the widespread application of such automated personality predictions from digital footprints among scientists and practitioners was hindered by the need to collect large amounts of self-report surveys in combination with textual data (see e.g., the myPersonality dataset [27]) to train and validate the predictive models. With the ability to make similar inferences with models that are available to the broader public, LLMs could democratize access to cheap and scalable psychometric assessments.

While this democratization holds remarkable opportunities for scientific discovery and personalized services, it also introduces considerable ethical challenges. Specifically, the ability to predict people's intimate psychological needs and preferences without their knowledge or consent poses a threat to people's privacy and self-determination [23]. For instance, users often share information online without considering how this information can be used by third parties and the use of LLMs for psychological profiling may not align with their original intentions. As the case of Cambridge Analytica [49] alongside a growing body of research on personalized persuasion and psychological targeting [47, 50, 51] has highlighted, insights into people's psychological dispositions can easily be weaponized to sway opinions and change behavior. Consequently, it might be necessary to introduce guardrails into systems like LLMs that prevent actors from obtaining psychological profiles of thousands or millions of users. Notably, the outlined concerns are aligned with recent calls for regulation [24–26] and the fact that the EU AI Act [52] explicitly bans emotion recognition in the workplace and educational institutions, as well as social scoring based on social behavior or personal characteristics.

4.4 Conclusion

Taken together, our research demonstrates the capacity of LLMs to derive psychological profiles from social media data, even without specific training. This zero-shot capability underscores the remarkable advancement LLMs represent in the domain of text analysis. While this "intuitive" understanding mirrors distinctly human abilities, the mechanisms and inherent biases associated with LLM-based personality judgments remain elusive and warrant further research. From a practical perspective, the potential of LLMs to effectively infer psychological traits from digital footprints presents a shift in psychometric evaluations, paving the way for large-scale AI-driven assessments. The prospect of democratized, scalable psychometric tools will enable breakthroughs in personalized services and large-scale research. Nevertheless, these advancements bring forth ethical challenges. The potential for non-consensual psychological predictions and other misuses highlights the necessity for stringent ethical frameworks.

Large Language Models Can Infer Psychological Dispositions of Social Media Users

Acknowledgments

We thank the Digital Future Initiative and Columbia Business School for their generous support. We thank Michal Kosinski for fruitful conversations and advice.

Author Contributions

H.P.: Conceptualization, Methodology, Software, Formal analysis, Investigation, Writing - Original Draft, Visualization; S.C.M.: Conceptualization, Methodology, Writing - Original Draft, Visualization.

References

- [1] OpenAI, GPT-4 Technical Report, 2023. [Online]. Available: https://cdn.openai.com/papers/gpt-4.pdf (visited on 08/21/2023).
- [2] Anthropic, *Model Card and Evaluations for Claude Models*, 2023. [Online]. Available: https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf.
- [3] T. Brown et al., "Language Models are Few-Shot Learners," in Advances in Neural Information Processing Systems, vol. 33, 2020, pp. 1877–1901. [Online]. Available: https://papers.nips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html (visited on 01/25/2024).
- [4] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language Models are Unsupervised Multitask Learners," 2019. [Online]. Available: https://www.semanticscholar.org/paper/Language-Models-are-Unsupervised-Multitask-Learners-Radford-Wu/9405cc0d6169988371b2755e573cc28650d14dfe (visited on 08/21/2023).
- [5] M. Kosinski, "Theory of Mind Might Have Spontaneously Emerged in Large Language Models," arXiv preprint arXiv:2302.02083, 2023. [Online]. Available: http://arxiv.org/abs/2302.02083 (visited on 09/08/2023).
- [6] T. Hagendorff, S. Fabi, and M. Kosinski, "Thinking Fast and Slow in Large Language Models," *arXiv* preprint arXiv:2212.05206, 2023. [Online]. Available: http://arxiv.org/abs/2212.05206 (visited on 09/08/2023).
- [7] J. Digutsch and M. Kosinski, "Overlap in meaning is a stronger predictor of semantic activation in GPT-3 than in humans," *Scientific Reports*, vol. 13, no. 1, p. 5035, 2023. [Online]. Available: https://www.nature.com/articles/s41598-023-32248-6 (visited on 09/08/2023).
- [8] S. C. Matz, J. D. Teeny, S. S. Vaid, H. Peters, G. M. Harari, and M. Cerf, "The potential of generative AI for personalized persuasion at scale," *Scientific Reports*, vol. 14, no. 1, p. 4692, 2024, Publisher: Nature Publishing Group. [Online]. Available: https://www.nature.com/articles/s41598-024-53755-0 (visited on 03/01/2024).
- [9] L. Albright, D. A. Kenny, and T. E. Malloy, "Consensus in personality judgments at zero acquaintance," *Journal of Personality and Social Psychology*, vol. 55, no. 3, pp. 387–395, 1988.
- [10] D. A. Kenny, L. Albright, T. E. Malloy, and D. A. Kashy, "Consensus in interpersonal perception: Acquaintance and the big five," *Psychological Bulletin*, vol. 116, no. 2, pp. 245–258, 1994, Place: US Publisher: American Psychological Association.
- [11] S. D. Gosling, S. J. Ko, T. Mannarelli, and M. E. Morris, "A room with a cue: Personality judgments based on offices and bedrooms," *Journal of Personality and Social Psychology*, vol. 82, no. 3, pp. 379–398, 2002.
- [12] P. J. Rentfrow and S. D. Gosling, "Message in a Ballad: The Role of Music Preferences in Interpersonal Perception," *Psychological Science*, vol. 17, no. 3, pp. 236–242, 2006. [Online]. Available: https://doi.org/10.1111/j.1467-9280.2006.01691.x (visited on 09/08/2023).
- [13] M. D. Back et al., "Facebook Profiles Reflect Actual Personality, Not Self-Idealization," Psychological Science, 2010. [Online]. Available: https://journals.sagepub.com/doi/epub/10.1177/ 0956797609360756 (visited on 01/26/2024).
- [14] M. Kosinski, D. Stillwell, and T. Graepel, "Private traits and attributes are predictable from digital records of human behavior," *Proceedings of the National Academy of Sciences*, vol. 110, no. 15, pp. 5802–5805, 2013.
- [15] D. Azucar, D. Marengo, and M. Settanni, "Predicting the Big 5 personality traits from digital footprints on social media: A meta-analysis," *Personality and Individual Differences*, vol. 124, pp. 150–159, 2018.
- [16] H. A. Schwartz *et al.*, "Personality, Gender, and Age in the Language of Social Media: The Open-Vocabulary Approach," *PLOS ONE*, vol. 8, no. 9, e73791, 2013. [Online]. Available: https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0073791 (visited on 09/07/2023).
- [17] G. Park *et al.*, "Automatic personality assessment through social media language," *Journal of Personality and Social Psychology*, vol. 108, no. 6, pp. 934–952, 2015.
- [18] T. Yarkoni, "Personality in 100,000 Words: A large-scale analysis of personality and word use among bloggers," *Journal of research in personality*, vol. 44, no. 3, pp. 363–373, 2010. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2885844/ (visited on 09/07/2023).

- [19] E. Grunenberg, H. Peters, M. J. Francis, M. D. Back, and S. C. Matz, "Machine learning in recruiting: Predicting personality from CVs and short text responses," *Frontiers in Social Psychology*, vol. 1, 2024, Publisher: Frontiers. [Online]. Available: https://www.frontiersin.org/articles/10. 3389/frsps.2023.1290295 (visited on 05/17/2024).
- [20] R. R. McCrae and P. T. Costa Jr., "The five-factor theory of personality," in *Handbook of personality: Theory and research, 3rd ed,* New York, NY, US: The Guilford Press, 2008, pp. 159–181.
- [21] T. Bolukbasi, K.-W. Chang, J. Zou, V. Saligrama, and A. Kalai, "Man is to computer programmer as woman is to homemaker? debiasing word embeddings," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, pp. 4356–4364. (visited on 01/25/2024).
- [22] Y. Wan, W. Wang, P. He, J. Gu, H. Bai, and M. Lyu, "BiasAsker: Measuring the Bias in Conversational AI System," *arXiv preprint arXiv:2305.12434*, 2023. [Online]. Available: http://arxiv.org/abs/2305.12434 (visited on 09/05/2023).
- [23] S. C. Matz, R. E. Appel, and M. Kosinski, "Privacy in the age of psychological targeting," *Current Opinion in Psychology*, vol. 31, pp. 116–121, 2020.
- [24] M. Perc, M. Ozer, and J. Hojnik, "Social and juristic challenges of artificial intelligence," *Palgrave Communications*, vol. 5, no. 1, s41599-019-0278-x, 2019. [Online]. Available: https://www.nature.com/articles/s41599-019-0278-x (visited on 01/24/2024).
- [25] P. Hacker, A. Engel, and M. Mauer, "Regulating ChatGPT and other Large Generative AI Models," in *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, 2023, pp. 1112–1123. [Online]. Available: https://dl.acm.org/doi/10.1145/3593013.3594067 (visited on 01/24/2024).
- [26] A. Chan, "GPT-3 and InstructGPT: Technological dystopianism, utopianism, and "Contextual" perspectives in AI ethics and industry," *AI and Ethics*, vol. 3, no. 1, pp. 53–64, 2023. [Online]. Available: https://doi.org/10.1007/s43681-022-00148-6 (visited on 01/24/2024).
- [27] M. Kosinski, S. C. Matz, S. D. Gosling, V. Popov, and D. Stillwell, "Facebook as a research tool for the social sciences: Opportunities, challenges, ethical considerations, and practical guidelines," *The American Psychologist*, vol. 70, no. 6, pp. 543–556, 2015.
- [28] L. R. Goldberg *et al.*, "The international personality item pool and the future of public-domain personality measures," *Journal of Research in Personality*, vol. 40, no. 1, pp. 84–96, 2006. [Online]. Available: https://linkinghub.elsevier.com/retrieve/pii/S0092656605000553 (visited on 09/08/2023).
- [29] A. Feingold, "Gender differences in personality: A meta-analysis," *Psychological Bulletin*, vol. 116, no. 3, pp. 429–456, 1994.
- [30] P. T. Costa Jr., A. Terracciano, and R. R. McCrae, "Gender differences in personality traits across cultures: Robust and surprising findings," *Journal of Personality and Social Psychology*, vol. 81, no. 2, pp. 322–331, 2001.
- [31] W. Youyou, M. Kosinski, and D. Stillwell, "Computer-based personality judgments are more accurate than those made by humans," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 112, no. 4, pp. 1036–1040, 2015. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4313801/ (visited on 01/25/2022).
- [32] S. Vazire and S. D. Gosling, "E-Perceptions: Personality Impressions Based on Personal Websites," *Journal of Personality and Social Psychology*, vol. 87, no. 1, pp. 123–132, 2004.
- [33] S. Abdurahman *et al.*, "Perils and Opportunities in Using Large Language Models in Psychological Research," *https://osf.io/tg79n*, 2023. [Online]. Available: https://osf.io/tg79n (visited on 05/12/2024).
- [34] E. Durmus *et al.*, "Towards Measuring the Representation of Subjective Global Opinions in Language Models," *arXiv:2306.16388* [cs], 2024, arXiv:2306.16388 [cs]. [Online]. Available: http://arxiv.org/abs/2306.16388 (visited on 05/13/2024).
- [35] S. Santurkar, E. Durmus, F. Ladhak, C. Lee, P. Liang, and T. Hashimoto, "Whose Opinions Do Language Models Reflect?" *arXiv:2303.17548* [cs], 2023, arXiv:2303.17548 [cs]. [Online]. Available: http://arxiv.org/abs/2303.17548 (visited on 05/13/2024).
- [36] M. Atari, M. J. Xue, P. S. Park, D. Blasi, and J. Henrich, "Which Humans?" https://osf.io/5b26t, 2023. [Online]. Available: https://osf.io/5b26t (visited on 05/12/2024).
- [37] S. Rathje, D.-M. Mirea, I. Sucholutsky, R. Marjieh, C. Robertson, and J. J. V. Bavel, "GPT is an effective tool for multilingual psychological text analysis," https://osf.io/sekf5, 2023. [Online]. Available: https://osf.io/sekf5 (visited on 05/12/2024).

- [38] S. E. Thayer and S. Ray, "Online Communication Preferences across Age, Gender, and Duration of Internet Use," *CyberPsychology & Behavior*, vol. 9, no. 4, pp. 432–440, 2006, Publisher: Mary Ann Liebert, Inc., publishers. [Online]. Available: https://www.liebertpub.com/doi/abs/10. 1089/cpb.2006.9.432 (visited on 05/12/2024).
- [39] K. Kondakciu, M. Souto, and L. T. Zayer, "Self-presentation and gender on social media: An exploration of the expression of "authentic selves"," *Qualitative Market Research: An International Journal*, vol. 25, no. 1, pp. 80–99, 2021, Publisher: Emerald Publishing Limited. [Online]. Available: https://doi.org/10.1108/QMR-03-2021-0039 (visited on 05/12/2024).
- [40] S. Tifferet and I. Vilnai-Yavetz, "Gender differences in Facebook self-presentation: An international randomized study," *Computers in Human Behavior*, vol. 35, pp. 388–399, 2014. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0747563214001381 (visited on 05/12/2024).
- [41] U. Oberst, V. Renau, A. Chamarro, and X. Carbonell, "Gender stereotypes in Facebook profiles: Are women more female online?" Computers in Human Behavior, vol. 60, pp. 559–564, 2016. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0747563216301480 (visited on 05/12/2024).
- [42] G. Roberti, "Female influencers: Analyzing the social media representation of female subjectivity in Italy," *Frontiers in Sociology*, vol. 7, 2022, Publisher: Frontiers. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fsoc.2022.1024043 (visited on 05/12/2024).
- [43] H. Peters, M. Cerf, and S. Matz, "Large Language Models Can Infer Personality from Free-Form User Interactions," 2024, Publisher: OSF Preprints. [Online]. Available: https://osf.io/apc5g/(visited on 05/19/2024).
- [44] T. Yang *et al.*, "PsyCoT: Psychological Questionnaire as Powerful Chain-of-Thought for Personality Detection," *arXiv:2310.20256* [cs], 2023, arXiv:2310.20256 [cs]. [Online]. Available: http://arxiv.org/abs/2310.20256 (visited on 04/17/2024).
- [45] S. R. Karra, S. T. Nguyen, and T. Tulabandhula, "Estimating the Personality of White-Box Language Models," arXiv:2204.12000 [cs], 2023, arXiv:2204.12000 [cs]. [Online]. Available: http://arxiv.org/abs/2204.12000 (visited on 04/23/2024).
- [46] P. M. Podsakoff, S. B. MacKenzie, and N. P. Podsakoff, "Sources of Method Bias in Social Science Research and Recommendations on How to Control It," *Annual Review of Psychology*, vol. 63, no. 1, pp. 539–569, 2012. [Online]. Available: https://doi.org/10.1146/annurev-psych-120710-100452 (visited on 09/07/2023).
- [47] S. C. Matz, M. Kosinski, G. Nave, and D. J. Stillwell, "Psychological targeting as an effective approach to digital mass persuasion," *Proceedings of the National Academy of Sciences*, vol. 114, no. 48, pp. 12714–12719, 2017. [Online]. Available: https://www.pnas.org/doi/10.1073/pnas.1710966114 (visited on 09/07/2023).
- [48] B. Freiberg and S. C. Matz, "Founder personality and entrepreneurial outcomes: A large-scale field study of technology startups," *Proceedings of the National Academy of Sciences*, vol. 120, no. 19, e2215829120, 2023. [Online]. Available: https://www.pnas.org/doi/10.1073/pnas. 2215829120 (visited on 09/08/2023).
- [49] M. Hu, "Cambridge Analytica's black box," Big Data & Society, vol. 7, no. 2, p. 2053 951 720 938 091, 2020. [Online]. Available: https://doi.org/10.1177/2053951720938091 (visited on 09/07/2023).
- [50] J. D. Teeny, J. J. Siev, P. Briñol, and R. E. Petty, "A Review and Conceptual Framework for Understanding Personalized Matching Effects in Persuasion," *Journal of Consumer Psychology*, vol. 31, no. 2, pp. 382–414, 2021. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/jcpy.1198 (visited on 09/07/2023).
- [51] M. Feinberg and R. Willer, "Moral reframing: A technique for effective and persuasive communication across political divides," *Social and Personality Psychology Compass*, vol. 13, no. 12, e12501, 2019.
- [52] E. Parliament, Artificial Intelligence Act: Deal on comprehensive rules for trustworthy AI, 2023. [Online]. Available: https://www.europarl.europa.eu/news/en/press-room/20231206IPR15699/artificial-intelligence-act-deal-on-comprehensive-rules-for-trustworthy-ai (visited on 01/25/2024).