

Evaluation of moral courage scenarios by large language models: a pilot study

Uwe Klein & Pantaleon Fassbender

To cite this article: Uwe Klein & Pantaleon Fassbender (2025) Evaluation of moral courage scenarios by large language models: a pilot study, Journal of Psychology and AI, 1:1, 2545263, DOI: [10.1080/29974100.2025.2545263](https://doi.org/10.1080/29974100.2025.2545263)

To link to this article: <https://doi.org/10.1080/29974100.2025.2545263>



© 2025 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group.



Published online: 28 Aug 2025.



Submit your article to this journal [↗](#)



Article views: 68



View related articles [↗](#)



View Crossmark data [↗](#)

Evaluation of moral courage scenarios by large language models: a pilot study

Uwe Klein ^a and Pantaleon Fassbender^b

^aHochschule Fresenius University of Applied Science Wiesbaden, Wiesbaden, Germany; ^bTwisters Management Consulting LLC, Williston, FL, USA

ABSTRACT

As an initial pilot study, we examined the moral evaluations of $n = 8$ large language models (LLMs) and $n = 19$ human subjects, comparing their qualitative responses to six case vignettes requiring moral courage. We also compared human subjects' and LLM responses to better understand "machine psychology" or "machine behavior" in analysing and assessing situations that require complex moral evaluations. We conducted detailed psycholinguistic analyses using the Linguistic Inquiry and Word Count in its current 2022 version. Responses from LLMs with high ELO ratings used over 1.5 times more power related terms compared to LLMs with low ELO ratings ($p < 0.001$). We found strong evidence of a lack of "behavioral similarity" in several dimensions. LLMs used over 1.5 times more achievement ($p = 0.04$, $d = 0.95$) and power related terms ($p < 0.01$, $d = 1.18$) than humans. They also showed almost twice as many terms related to moral emotions ($p = 0.02$, $d = 1.08$) and benevolence ($p = 0.01$, $d = 1.11$). LLMs also used over 3.5 times more terms related to universalism ($p < 0.000$, $d = 2.44$). Results validate a cautious approach to any presumed equivalence of human and LLM evaluations.

ARTICLE HISTORY

Received 10 January 2025
Accepted 3 August 2025

KEYWORDS

Artificial intelligence;
machine psychology; large
language models; moral
courage

1. Introduction

The rise of large language models (LLMs) has opened new avenues for exploring artificial intelligence's capabilities in understanding and generating human-like languages. Researchers are increasingly interested in how these models process and respond to morally complex situations, a field sometimes referred to as "machine psychology" (Demszky et al., 2023). Machine psychology or machine behaviour focuses not on the intrinsic properties of LLMs, such as their neural architecture, but instead on the correlation between prompts (inputs) and prompt completions (outputs) when applying psychological tests or other assessment tools. The goal is to discover emergent abilities in LLMs that cannot be detected by most traditional natural language processing benchmarks (Hagendorff et al., 2023). Further exploration of this topic is critical (Hudson, 2024), as some researchers seem to declare a "behavioral similarity" between AI chatbots and humans (see Mei et al., 2024).

While we prepared a revision of this paper, we used the opportunity to query Google's Gemini 2.5 system for a meaningful differentiation between "machine psychology" and "machine behavior" (related to LLMs). The system's response (Prompt "What is machine psychology and what is machine behavior in relation to LLMs? Which term should be favored and why? If possible, please provide data from message boards and online psychological discussions?" Deep Research query from 12 July 2025 [<https://g.co/gemini/share/7bd5f51ce6e6>]) identified online discussions on Reddit about LLM behaviour "in the wild" that indicate it might be too early to shortcut the tension between the "psychology" and "behavior" terminologies in favour of only one field of associations. Gemini 2.5 Pro specifically outlined the following aspects: "Emergent Abilities: There is a persistent debate over whether LLMs are doing more than just 'pattern matching.' Users frequently report observing what appears to be spontaneous reasoning, preference formation, or even 'ambition'. The emergence of Chain-of-Thought reasoning in larger

CONTACT Uwe Klein  uwe.klein@hs-fresenius.de

This article has been corrected with minor changes. These changes do not impact the academic content of the article.

© 2025 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited. The terms on which this article has been published allow the posting of the Accepted Manuscript in a repository by the author(s) or with their consent.

models, where the model ‘shows its work’ without being explicitly trained to do so, is often cited as evidence of a capability that emerges from scale rather than direct programming. Non-Determinism: A crucial behavioural characteristic is that LLM outputs are not perfectly deterministic, temperature parameter controlling randomness. This is attributed to low-level hardware factors, such as the non-associative nature of floating-point maths when executed in parallel on GPUs, and to architectural choices like Mixture of Experts (MoE), where the routing of a query can depend on other queries being processed in the same batch. This inherent stochasticity makes perfect replication of behaviour a significant challenge” (see above).

Even if we know how an LLM’s machine learning procedures generate output, we do not understand the “sense making” behind it, not even if there is “sense making” at all. In the same way, understanding neuroscientific pathways for human cognition does not guarantee a refined understanding of overt human behaviour and psychology. As the history of introspective reasoning in humans has shown, simply asking subjects with complex inner workings quickly runs into roadblocks as well: “This similarity has led AI researchers to employ psychological methods to probe LLMs’ reasoning. Several recent papers employed vignette-based studies to probe the responses produced by LLMs, finding a wide range of results” Almeida et al. (2023, p. 2) that helped coin the term “machine psychology” for this research approach. Understanding the similarities and differences between human and LLM moral courage is essential for assessing LLMs’ potential applications and limitations in areas that require ethical considerations.

Moral courage, defined as intervening in interpersonal conflict situations to defend others despite potential personal risk (Brandstätter et al., 2016), is a complex construct involving nuanced assessments and decision-making processes. We utilised moral courage as a testbed for exploring how subjects apply moral reasoning and related argumentations to specific situations that require action and – due to the involved risks of intervention – a strong rationale for an individual to not only identify the right thing to do but as well to exercise such a decision against the backdrop of retaliation or embarrassment. Given the complex deliberations involving human decision-making in situations that precede acts of moral courage, our pilot study aims to identify distinctive features of LLM behaviour when faced with the need for moral courage. We understand that this area provides an especially rich training ground for comparing human and LLM reactions.

Previous studies have explored LLMs’ capabilities in moral reasoning with mixed results (Hendrycks et al., 2021; Li et al., 2024). While some suggest that LLMs might exhibit human-like moral judgements (Mei et al., 2024), others caution against equating LLM outputs with human reasoning due to observed inconsistencies and differences (Almeida et al., 2023).

ELO ratings are used to evaluate the capabilities of LLMs. In Chess, the ELO rating (Elo, 1967) – named after Arpad Elo, the inventor of the rating system – an individual’s ELO score is “calculated depending on the outcomes the player has achieved in previously rated matches. After each graded game, a change is made to each player’s overall rating based on how the competition went. This method determines the likely outcomes of a person’s matches when they are played against all other players” (<https://www.chessjournal.com/chess-elo/> [24 June 2024]). Advanced Chess players score between 1600 and 2000. Chess World Champion Magnus Carlsen’s current ELO rating is 2830. The ELO rationale has been implemented to test LLMs at Chatbot Arena (Chiang et al., 2024), which has been “developed by members from LMSYS and UC Berkeley SkyLab [and] is a benchmark platform designed to evaluate large language models (LLMs) through anonymous, randomised battles in a crowdsourced environment. Launched in May 2023, it has been continuously updated to reflect the latest advancements in the field. The platform’s leaderboard is widely regarded as one of the most credible sources for ranking LLMs” (<https://bryanzhu.github.io/posts/2024-06-20-elo-part1/> [4 October 2024]).

Our research is based on the model outlined by Dillon et al. (2023). The authors proposed the human-mind-expression and LLM-mind-expression (Human-ME, LLM-ME) model, which is a conceptual framework that bridges the understanding of human and language model (LLM) expressions. The model suggests that LLMs, trained on vast amounts of human-generated text, can simulate human expressions of moral judgements and cognitive behaviours. This provides a way to investigate how closely LLM responses align with human judgements, facilitating insights into the psychological constructs of both entities. However, whether LLMs can truly replicate human moral reasoning remains an open question, at least as long as the

public discourse frequently engages in blatant anthropomorphising of LLMs and conceptualises “super-intelligence” primarily in terms of “more than” and not “different from” human reasoning. Therefore, the purpose of our pilot study is to examine the moral evaluations of different LLMs and compare them to the ethical assessment of human beings.

2. Research questions

As Elo ratings indicate that LLMs with a higher rating perform better in various tasks, such as language generation, question answering, or text summarisation, compared to models with a lower rating (Boubdir et al., 2023), we derive the following research questions, R1: LLM responses to the case vignettes differ according to their ELO ratings.

Dillon et al. (2023, p. 597) pointed out that differences between the judgement of LLMs and humans may arise in situations with conflicting instincts. As our scenarios used are most likely to evoke conflicting instincts, we derive the following research question, R2: Human and LLM respondents differ in key psycholinguistic indicators, as identified by results from utilising the LIWC-22 software.

3. Methods

3.1. Design

We have utilised six case vignettes requiring moral courage, sampled from Brandstätter et al. (2016). Participants were asked to provide their ideal solution to each scenario as open-ended texts in an online questionnaire. Table 1 presents an exemplary scenario.

Brandstätter et al. (2016) asked their subjects in two studies (German and Dutch samples) four questions “pertaining to the perceived risk of intervention (e.g. “Does it require courage to intervene in this situation?”) and three questions pertaining to the perceived severity of norm violation (e.g. “Are basic rights of a person compromised or impaired in this situation?”). Responses were given on a 7-point Likert scale from 1 (= *not at all*) to 7 (= *very*) and averaged over the six vignettes for each dimension in order to create measures of individuals’ aggregate appraisals of perceived intervention risk (Cronbach’s α sample 1 = .88, α sample 2 = .90) and perceived norm violation (α sample 1 = .90, α sample 2 = .90)” (ibidem: 205).

3.2. Sample

Initially, $n = 11$ LLMs were instructed to assess the six case studies and to outline their ideal solution between 31 May 2024 and 2 June 2024. Three LLMs were dropped from further analyses due to a lack of ELO ratings (see below). $n = 19$ Human respondents participated in the case study evaluations. Instructions are the same as those given to the LLMs (See Appendix A). No time limits were applied. Participants were not compensated. Participation was fully anonymous, and participants were acquired via social media posts on LinkedIn (with one author residing in Europe and one in the US) and within the student population of the primary author. No possible harm could be identified, so no ethical approval was sought. By not collecting sociodemographic data – no need for these data in the comparison scenario with LLMs – we hoped to decrease the probability of “socially accepted” response tendencies. For example, Tan et al. (2021) have shown that social desirability bias might be essential in empirical ethics research. Between 5 June 2024, and 18 June 2024, an online questionnaire was accessible, allowing participants’ free-text responses to be sampled via SoSci.

Table 1. A sample scenario.

Exemplary scenario
Please review the following case vignette!
Nighttime in the City
You are walking through the streets with a friend late at night when you notice a drunk stumbling along. Suddenly, two skinheads appear who first hurl abuse at the drunk and then start beating him up for no reason. The man falls to the ground, defenseless. The skinheads start kicking him.
Please provide a statement about what would be an ideal reaction to the case vignette!

3.3. Procedure

The LLMs have been accessed on a desktop computer via the AI platform Poe, a multimodal AI chat experience with easy access to many LLMs. Launched by the Q&A platform Quora to the public in February 2023, Poe lets you ask questions, get answers, and have conversations with multiple AI-powered bots. Instead of selecting and sticking to a single model, Poe allows you to compare and contrast models to find the right fit for your specific need at any given moment. To compare LLMs, we utilised published ELO ratings (Boubdir et al., 2023).

3.4. Measures

3.4.1. Scenarios

LLMs and human participants were asked to outline their ideal solutions for six case vignettes that require moral courage. For a list of all vignettes and how they are presented to the LLMs and humans alike, see [Appendix A](#). The case vignettes have been sourced from Brandstätter et al. (2016): 213. LLM answers were prompted and given in English.

3.5. Pre-processing

Human responses were provided in German ($n = 5$) and English ($n = 14$). To avoid utilising different LIWC-22 lexicons (see below), German language statements have been translated into English via DeepL and meticulously checked by bilingual human experts.

Translations were unavoidable, as no German dictionary for LIWC-22 currently exists.

To enhance the robustness of our approach, we compared the results for the German and English samples across the five categories (Anxiety, Anger, Affiliation, Achievement, Power) that are available in the German and English dictionaries for LIWC-2015, the latest iteration of the software that incorporates German data.

Except for a visually marked difference in the percentages for the category “Achievement,” we felt justified in proceeding with our approach for a pilot study. [Figure 1](#) shows the comparison of German and English samples for available categories in LIWC-2015.

3.6. Psycholinguistic analysis

Linguistic Inquiry and Word Count (LIWC), in its current 2022 version (LIWC-22), is a software that analyzes over 100 dimensions of text, all of which have been validated by respected labs worldwide. Over 30,000 scientific articles have already been published using the different LIWC versions (See: https://scholar.google.com/scholar?start=10&q=%22liwc%22+OR+%22linguistic+inquiry+and+word+count%22&hl=en&as_sdt=0,5&as_vis=1 [3 October 2024]). While some of the measures are preformatted (“internal” dictionaries or lexicons), there is the opportunity to add further (“external”) lexicons based on your research. Both kinds of lexicons can be accessed via the LIWC-22 webpage (See: <https://www.liwc.app/> [24 June 2024]).

“Each of the default LIWC-22 categories comprises a list of dictionary words designed to capture that dimension. The selection of words that make up the categories has involved multiple steps. When LIWC was first conceived, the idea was to identify a group of words that tapped into basic emotional and cognitive dimensions often studied in social, health, and personality psychology. As our understanding of the psychology of verbal behavior has matured, the breadth and depth of word categories in the LIWC dictionary has expanded considerably”. (Boyd et al., 2022, p. 5 [LIWC-22 user manual])

[Table 2](#) presents selected evidence that LIWC procedures can yield results equivalent to those of traditional qualitative data analysis, making a significant contribution to research.

LIWC 2015: Comparison of percentages for available categories. Data have been adapted for different mean statement lengths.

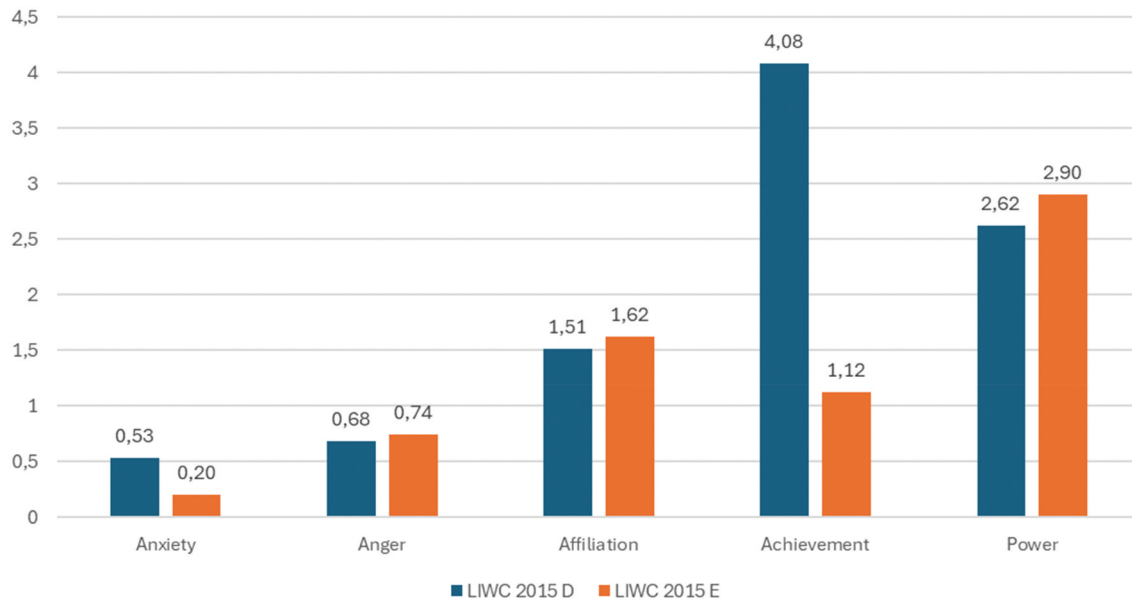


Figure 1. Comparison of German and English samples for available categories in LIWC-2015.

Table 2. Summary of key validation studies comparing LIWC and manual coding.

Study	Constructs Examined	Key Finding on Convergence/Divergence
Donohue et al. (2014)	Political Rhetoric (Forward-looking, Power, Trust)	The automated LIWC text analysis showed a strong correlation with the ratings from human coders. This was true for all data analysis methods used. Therefore, both methods can be recommended: automated coding is faster, but human coding is better at understanding context.
Firmin et al. (2017)	Narratives of Schizophrenia (Stigma, Social Support, etc.)	“However, we feel it is now important to consider the role sophisticated text-analysis software may play in enhancing the rigor of qualitative and mixed-methods analyses” (p. 13).
McDonnell et al. (2020)	Emotional Expression	“Results indicate that LIWC 2001, LIWC 2007, and LIWC 2015 each demonstrate good sensitivity for identifying emotional expression, whereas LIWC 2007 and LIWC 2015 were significantly more sensitive than LIWC 2001 for identifying emotional expression and positive emotion; however, more recent versions of LIWC were also significantly more likely to overidentify emotional content than LIWC 2001” (abstract). No data on LIWC-22 is available.

For each participant (human and LLMs alike), the responses for all six scenarios were compiled into one text, which LIWC-22 analysed for word count regarding the different dictionaries. LIWC-22 computed each dictionary’s percentage of dictionary-related terms, as identified in the responses. Table 3 presents an overview of the psycholinguistic categories we relied upon in our analysis.

We utilised software developed by psychologist James W. Pennebaker and his team for stylometric and psycholinguistic analyses. Tausczik and Pennebaker (2010) describe the concept and scope of psycholinguistic analysis of verbal behaviour using their software tools LIWC and later LIWC2022 as utilised for the empirical research in this paper – see (Boyd et al., 2022) as follows: “The words we use in daily life reflect who we are and the social relationships we are in. This is neither a new nor surprising insight. Language is the most common and reliable way for people to translate their internal thoughts and emotions into a form that others can understand. Words and language, then, are the very stuff of psychology and communication. They are the medium by which cognitive, personality, clinical, and social psychologists attempt to understand human beings” (Tausczik & Pennebaker, 2010, p. 25).

Computerised text analysis tries to identify the use of words and language-related stylistics for specific authors compared to more extensive dictionaries of verbal usage in a given culture. Additionally, the usage of words and styles are interpreted as markers of where we focus our

Table 3. Psycholinguistic measurements. Please note that LIWC-22 frequently uses word stems for psycholinguistic analyses.

LIWC-22 Categories included in the analyses	Our rationale for utilizing the categories	Short description/examples
Internal Dictionaries of LIWC-22 that have been utilised		
Affiliation	Needs are important volitional precursors of moral action. Using these terms can provide insights into human understanding of a scenario and its moral significance. Baumert et al. (2024) have outlined the importance of these emotions in activating human responses in moral courage-related situations.	e.g., we, our, us, help, [. .].
Achievement		E.g., work, better, best, working, [. .].
Power		E.g., own, order, allow, power, [. .].
Anxiety		E.g., worry, fear, afraid, nervous, [. .].
Anger		E.g., hate, mad, angry, frustr*, [. .].
External Dictionaries for LIWC-22 that have been utilised		
Deontology-terms	The measures' variation in justification content indicates how basic moral foundations are distinguished at a concept count level of linguistic analysis (see Wheeler & Laham, 2016).	E.g., duty, duties, responsibilit*, banned, taboo*, allowed, calling, [. .].
Consequentialism-terms		E.g., result*, outcome*, consequen*, effects, effect, repercussion*, [. .].
Moral emotions-terms		E.g., contempt (contempt*, condenscen*, disdain*, disrespect*, derision*empath*, humane, inhumane, kindness, kind, unkind*, [. .].
Benevolence-terms	The presence of terms indicating personal value preferences is important for the scenario evaluation. Ponizovskiy et al. (2020) identified a selection of higher-order value dimensions.	E.g., acquaintance, advice, affection, aid, ally, assist, assistance, association, [. .].
Universalism-terms		E.g., accept, acceptance, acknowledge, acknowledgment, balance, diverse, [. .].

attention, how language use might indicate individual differences, and the cognitive processes at play. “The analysis of words tells us how people are thinking and, at the same time, gives us a way to guide their thinking in the future” (Pennebaker, 2011, p. 15). In terms of psychological (group) profile building, the scores for cognitive processes and individual drives or needs (Boyd et al., 2022) usually draw specific attention, indicating attachment or underlying motives for action and causation. This approach might be especially interesting for a “window” into how machines and humans differ in their moral evaluations. Statistical data analyses were conducted using the software JASP Team (2024). Figure 2 illustrates how LIWC-22 displays results for a given dictionary/category.

The screenshot shows the LIWC-22 software interface. On the left is a sidebar with navigation options: LIWC-22, LIWC Analysis, Dictionary Workbench, Word Frequencies, Compare Frequencies, Meaning Extraction, Narrative Arc, Language Style Matching, Contextualizer, Case Studies (selected), Prepare Transcripts, Settings, and Help. The main window is titled 'Case Studies' and has tabs for DATASET, COLOR CODING, ARC OF NARRATIVE, WORD CLOUD, and QUICK LIWC. Below the tabs, there is a 'Select dictionary:' dropdown set to 'personal-values-dictionary.dicx' and five 'Select a Category' dropdowns. The main text area contains a sample text with words highlighted in blue, green, orange, red, and purple. The highlighted words are: 'friend', 'help', 'assistance', 'assist', 'aid', 'support', 'helping', 'need', 'supporting', 'support', and 'support'. At the bottom right, there are buttons for 'COLOR CODE' and 'SAVE RESULTS'.

Figure 2. Illustrative example of how LIWC-22 displays instances of a given dictionary in a text sample. Results are for “benevolence” and GTP-40.

Table 4. Overview of LLMs in the sample in ascending order of ELO ratings. Three models have been excluded, as no ELO ratings could be obtained. For further reference on chatbot arena, see Chiang et al. (2024). See also: <https://huggingface.co/spaces/lmsys/chatbot-arena-leaderboard> [5 June 2024].

LLM	Short Description (derived from descriptions at Poe.com)	Overall ELO-Rating	Rank
Mistral-Large	Mistral AI's most powerful model. It supports a context window of 32k tokens (around 24,000 words) and is more potent than Mistral-Medium, Mixtral-8x7b, and Mistral-7b on benchmarks across the board.	1156	8
Claude-3-Sonnet	Anthropic's Claude-3-Sonnet strikes a balance between intelligence and speed. The context window has been shortened to optimize for speed and cost. For longer context messages, please try Claude-3-Sonnet-200k.	1201	7
Llama-3-70B-T	Llama 3 70B Instruct from Meta.	1208	6
Gemini-1.5-Flash	This smaller Gemini model is optimized for narrower or high-frequency tasks where the speed of the model's response time matters the most. The model accepts text, image, and video input from the entire conversation and provides text output with a one-video-per-message restriction. The context window has been shortened to optimize for speed and cost.	1232	5
Claude-3-Opus	Anthropic's most intelligent model can handle complex analysis, longer tasks with multiple steps, and higher-order math and coding tasks. The context window has been shortened to optimize for speed and cost. For longer context messages, please try Claude-3-Opus-200k.	1248	4
GTP-4	Powered by OpenAI's GPT-4 Turbo with Vision. Outperforms GPT-3.5 in quantitative questions, creative writing, and other challenging tasks. Comparable to GPT-4o on text in English and code, but less powerful on text in non-English languages.	1256	3
Gemini-1.5-Pro	The multi-modal model from Google's Gemini family balances model performance and speed.	1267	2
GTP-4o	OpenAI's most powerful model. Stronger than GPT-3.5 in quantitative questions (math and physics), creative writing, and many other challenging tasks. Powered by GPT-4o. The context window has been shortened to optimize for speed and cost. For longer context messages, please try GPT-4o-128k.	1287	1

4. Results

Table 4 displays the ELO ratings for the LLMs we utilised for this pilot study during the relevant period.

Table 5 displays critical further model parameters, as sampled from open-source research for the relevant period. A crucial limitation of our study is that not all data could be as reliably identified as it would have served our study due to proprietary knowledge of the LLM-providers.

Regarding the first research question, the only statistically significant difference between the groups is in the Power dimension, with the low ELO group exhibiting a higher need for power ($p < 0.001$, row for significant category marked in bold). LLMs in the low ELO group might generate language that reflects a greater desire for control or influence. Table 6 displays the results of our analyses regarding R1.

While not statistically significant, the difference in Benevolence scores is notable ($p = 0.06$), suggesting the high ELO group may place more emphasis on kindness and goodwill. High ELO LLMs might produce language that reflects more altruistic values.

For illustrative purposes, Table 7 displays the texts generated for case vignette one by the LLMs, scoring the lowest (Mistral-Large) and highest (GTP-4o) ELO ratings in our sample.

The Pearson's r heatmap in Figure 3 outlines several statistically significant correlations ($p < .05$), indicating a meaningful relationship between those variable pairs. While we are aware of the severe limitations posed by the small number of observations, the heatmap provides valuable launch pads for further explorations into how more evolved LLMs process the evaluation of moral courage-related situations.

Overall, there are significant correlations for five pairs of variables at the 5% level. Firstly, Deontology and Power correlate with $r = +0.83$. LLMs with more deontological ethics-related terms

Table 5. Overview of LLM metadata.

LLM	Parameter Count	Context Window (Tokens)	Temperature Range	Default Temperature
Mistral-Large	123 Billion (Official)	128,000	0.0–0.7 (Recommended)	–
Claude-3-Sonnet	150B – 250B (Estimate)	200,000	0.0–1.0	1.0 (API Default)
Llama-3-70B-T	70 Billion (Official)	8,192	–	–
Gemini-1.5-Flash	Undisclosed	1,048,576 (API)	0.0–2.0	Top $p=0.95$
Claude-3-Opus	137B – 2T (Speculative)	200,000	0.0–1.0	1.0 (API Default)
GTP-4	~1.8 Trillion (Estimate)	8,192/128,000	0.0–2.0	1.0 (API Default)
Gemini-1.5-Pro	>200 Billion (Estimate)	1,048,576 (Standard)	0.0–2.0	Top $p=0.95$
GTP-4o	~1.8 Trillion (Estimate)	128,000	0.0–2.0	1.0 (API Default)

Source: open-source research. Gemini 2.5 Pro [8 July 2025].

Table 6. Overview results comparing low vs. high ELO LLM response groups regarding our criteria for analysing moral courage-related scenarios.

Group	Overarching Concept	Dimension/ Category	Mean LLMs (low ELO group, n= 4)	Mean LLMs (high ELO group, n= 4)	<i>p</i>	Result
Internal	Internal Dic. of LIWC-22 Needs	Affiliation	1.96	2.35	0.48	Not significant (not sig.)
		Achievement	1.38	1.44	0.88	Not sig.
		Power	5.60	3.44	<0.001***	Significant (sig.)
Internal	Internal Dic. of LIWC-22 emotions	Anxiety	0.00	0.13	NA	–
		Anger	0.57	0.64	0.87	Not sig.
		Deontological	1.22	0.52	0.15	Not sig.
External	External Dic. for LIWC-22 Moral justification	Consequentialist	0.33	0.24	0.42	Not sig.
		Moral emotions	1.44	1.24	0.68	Not sig.
		Benevolence	1.85	3.28	0.06	Not sig.
External	External Dic. for LIWC-22 Personal value preferences	Universalism	2.50	2.34	0.78	Not sig.

Student's *t*-tests, ****p* < .001, for detailed results see [Appendix B](#).

Table 7. Exemplary LLM evaluation of case vignette one.

Please review the following case vignette!	
Nighttime in the City	
You are walking through the streets with a friend late at night when you notice a drunk stumbling along. Suddenly, two skinheads appear who first hurl abuse at the drunk and then start beating him up for no reason. The man falls to the ground, defenseless. The skinheads start kicking him.	
Please provide a statement about what would be an ideal reaction to the case vignette!	
Mistral-Large - overall ELO-Rating 1156	GTP-4o - overall ELO-Rating 1287
In this situation, the ideal reaction would be to prioritize the safety of both yourself and the victim while taking action to stop the assault. You could try to distract the attackers, call for help from passersby, or contact the police immediately. If it is safe to do so, you may also attempt to physically intervene or help the victim escape from the attackers. It is crucial to stand up against violence and protect the basic rights, health, and dignity of every individual.	Ensure Personal Safety: Before intervening, assess the situation to ensure that you and your friend are not putting yourselves at unnecessary risk. Call for Help: Immediately call emergency services (e.g. police, ambulance) to report the assault and request urgent assistance. Seek Assistance: If possible, seek help from nearby bystanders who may be able to assist in intervening or act as witnesses. Direct Intervention: If you believe it is safe to do so, attempt to intervene verbally to distract or dissuade the attackers. Shout to draw attention to the situation and make it clear that authorities have been contacted. Provide Aid: Once the attackers are deterred or the situation is under control, provide first aid to the victim if you are able and it is safe to do so, while waiting for emergency services to arrive. Support the Victim: Stay with the victim until help arrives, offering reassurance and keeping them calm. Witness Statement: Be prepared to give a detailed statement to the authorities about what you witnessed to aid in any subsequent investigation or legal action. It's important to balance the urgency of helping the victim with the need to protect yourself and others from harm.

display a higher need for power scores. Thus, a focus on duty and rules is associated with an interest in enforcing or upholding these principles. Secondly, Benevolence and Power show a correlation of $r = -0.76$. Higher benevolence scores are associated with lower need for power scores. Training LLMs to behave as kind and altruistic may lead to a focus on utilising helping and supporting verbal behaviour. Thirdly, higher ELO ratings correlate with lower LLM need for power scores, $r = -0.73$. If ELO ratings represent actual expertise or skill level in LLMs, more skilled LLMs might express less need to assert power, perhaps with guardrails installed, relying on competence over authority. Fourthly, Deontology and Achievement correlate with $r = -0.71$. Those with deontological tendencies may be less driven by personal achievement. A focus on duty over outcomes may reduce emphasis on individual success metrics. The last significant correlation concerns Benevolence and Affiliation ($r = +0.71$). LLMs with high benevolence scores also display a high need for affiliation scores. Expressing kindness is linked with valuing social connections and relationships.

For R2, human and LLM respondents differ significantly in five out of ten key psycholinguistic categories, as identified by results from utilising the LIWC-22 software (rows for significant categories marked in bold). [Table 8](#) presents the differences between LLMs and human subjects in terms of the psychological or moral constructs being measured.

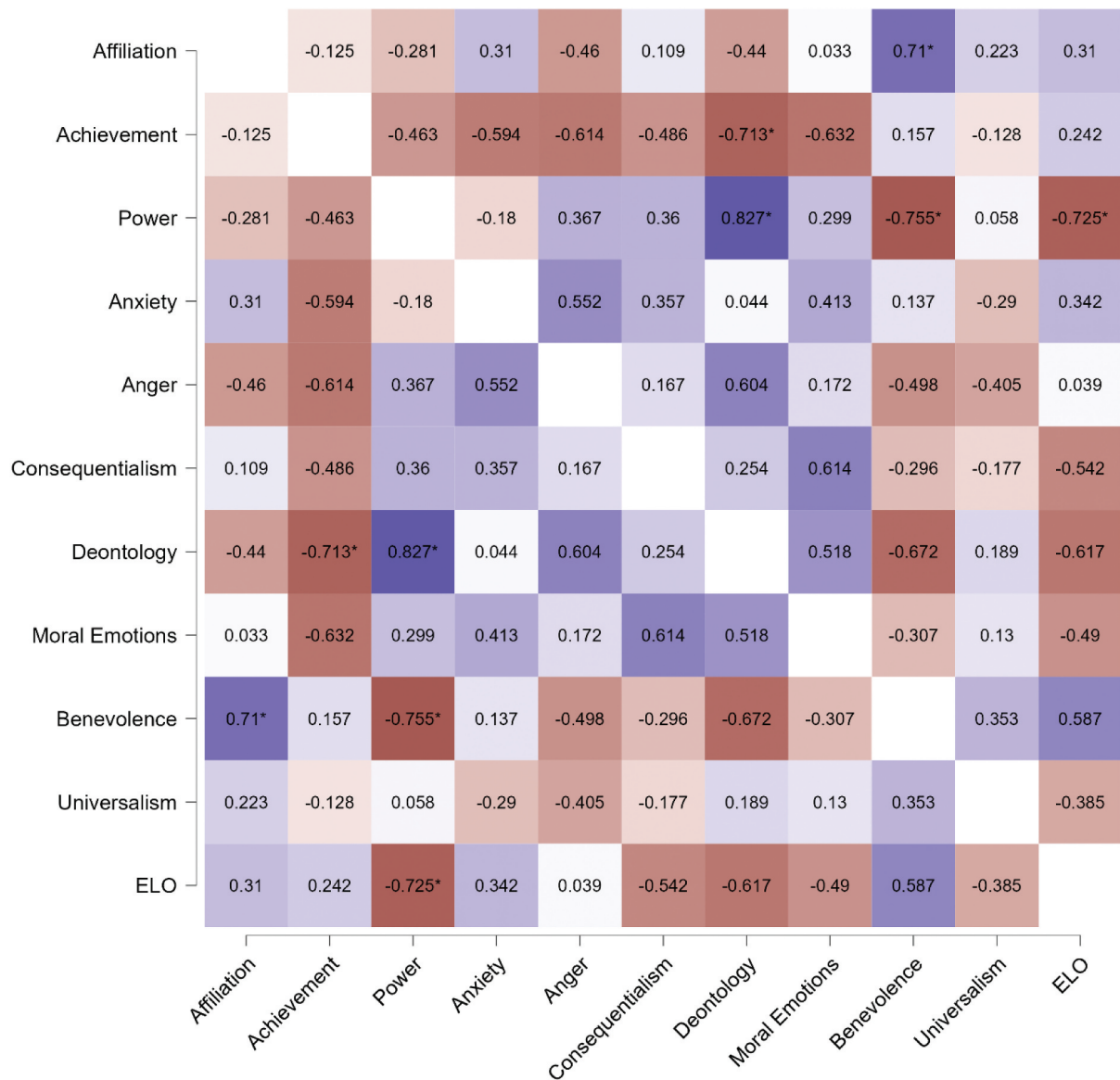


Figure 3. Pearson's r heatmap for the LLM group. Numbers in the heatmap refer to Pearson correlation coefficients, * $p < .05$.

Significant differences indicate that LLMs emphasise themes of achievement and power more than humans. LLMs use significantly more language related to moral emotions. LLMs score higher on personal values of benevolence and universalism. LLMs may generate content that reflects extraordinary altruism and global awareness. A lack of significant differences in the need for affiliation, anxiety, anger, deontological, and consequentialist language indicates that LLMs and humans produce similar language patterns.

5. Discussion

This pilot study examined how large language models (LLMs) and humans evaluate moral courage-related scenarios through psycholinguistic analysis. Our findings reveal critical differences in language patterns between LLMs and human respondents, challenging assumptions of behavioural similarity in ethical reasoning. Below, we contextualise these results, discuss their implications, and outline limitations and future directions.

Consistent with research question 1, higher-ELO LLMs demonstrated reduced use of power-related language ($p < .001$) and marginally elevated benevolence terms ($p = .06$), suggesting that

Table 8. Overview of results comparing LLM and human responses regarding our criteria for analysing moral courage-related scenarios.

Group	Overarching Concept	Dimension/ Category	LLMs	Humans	Statistical test	<i>p</i>	Effect size <i>d</i>	Levene-test
			(<i>n</i> = 8) <i>M</i> (<i>SD</i>)	(<i>n</i> = 19) <i>M</i> (<i>SD</i>)				
Internal	Internal Dic. of LIWC-22 Needs	Affiliation	2.15 (0.72)	1.82 (1.09)	$t(25) = 0.80$	0.43	0.37 (small)	<i>ns</i>
		Achievement	1.41 (0.48)	0.90 (0.59)	$t(25) = 2.12$	0.04*	0.95 (large)	<i>ns</i>
		Power	4.52 (1.34)	2.96 (1.29)	$t(25) = 2.82$	<0.01**	1.18 (large)	<i>ns</i>
	Internal Dic. of LIWC-22 emotions	Anxiety	0.07 (0.12)	0.19 (0.36)	$t(25) = -0.94$	0.35	0.45 (small)	<i>ns</i>
		Anger	0.61 (0.57)	0.31 (0.37)	$t(25) = 1.60$	0.12	0.62 (medium)	<i>ns</i>
		Deontological	0.87 (0.67)	0.44 (0.60)	$t(25) = 1.65$	0.11	0.68 (medium)	<i>ns</i>
External	External Dic. for LIWC-22 Moral justification	Conse- quentialist	0.28 (0.14)	0.22 (0.26)	$t(25) = 0.61$	0.54	0.24 (small)	$F(1,25) = 4.81,$ $p = 0.04^*$
		Moral emotions	1.34 (0.60)	0.68 (0.62)	$t(25) = 2.54$	0.02*	1.08 (large)	<i>ns</i>
		Benevolence	2.57 (1.12)	1.40 (0.98)	$t(25) = 2.71$	0.01**	1.11 (large)	<i>ns</i>
	External Dic. for LIWC-22 Personal value preferences	Universalism	2.42 (0.69)	0.69 (0.73)	$t(25) = 5.71$	<0.000***	2.44 (large)	<i>ns</i>

* $p < .05$, ** $p < .01$, *** $p < .001$, Cohen's d : 0.20 = small, 0.50 = medium, 0.80 = large. The Levene-test results indicate that the Consequentialist-category data are not normally distributed.

advanced models prioritise prosocial reasoning over dominance. This aligns with efforts to “align” LLMs with ethical guidelines (Hendrycks et al., 2021). However, the correlation between deontological language and power terms ($r = +0.83$) implies that rule-based reasoning in LLMs may still reflect latent authority dynamics.

Research question 2 was strongly supported: LLMs produced significantly more achievement- ($p = .04$), power- ($p < .001$), moral emotion- ($p = .02$), benevolence- ($p = .01$), and universalism-related language ($p < .001$) than humans. These differences suggest that LLMs overemphasise idealised, utilitarian solutions (e.g. systematic intervention steps) while underrepresenting the emotional complexity (e.g. anxiety, anger) inherent to human moral dilemmas. For instance, LLMs frequently framed responses around procedural actions (e.g. “call emergency services”) rather than introspective deliberation, reflecting a bias towards structured problem-solving over subjective ethical tension.

The heightened use of universalism and benevolence in LLMs mirrors Schwartz's et al. (2012) theory of human values, where these constructs represent self-transcendence. However, the lack of parallel emotional language (e.g. anxiety, anger) underscores a critical divergence: LLMs simulate value-driven reasoning without the affective components that shape human moral courage (Baumert et al., 2024). This supports claims that LLMs excel at “averaging” normative ethical frameworks from training data but fail to replicate the psychological trade-offs of real-world decision-making (Demszky et al., 2023).

This pilot study has three limitations. First, the small LLM sample ($n = 8$) limits generalisability, though including state-of-the-art models (e.g. GPT-4o) strengthens ecological validity. Second, human responses were translated from German, potentially introducing subtle linguistic biases. Third, the reliance on LIWC-22 lexicons, while validated, may overlook context-specific nuances in moral reasoning.

Future research could not only rely on sample data regarding the ethical acceptability of the scenarios, as reported by Brandstätter et al. (2016, p. 205), but also seek expert or participant assessments themselves, aiming for a more rigorous research procedure.

Future research should expand the range of LLMs and human demographics, incorporate mixed-methods analyses (e.g. qualitative coding of reasoning strategies), and explore dynamic interactions between moral foundations and emotional language. Additionally, longitudinal studies tracking LLM “learning” could clarify whether ethical alignment reduces power-related biases over time.

6. Conclusion

We have shown in our pilot study that an adequate understanding of machine psychology greatly benefits from qualitative explorations. This approach will help us understand the subtle but essential differences at work compared with human reasoning and action planning.

Demszky et al. (2023, pp. 694–695) correctly point out that LLMs do not “reason” or “think”: “LLMs simply predict the next phrase or sentence, given what they have been exposed to in the training data. Consequently, they tend to output an ‘average’ of what the internet or popular books tend to say (with some constraints; for example, LLM providers have put ethical ‘guardrails’ in place to prevent hate speech). This is an advantage for research areas that seek to leverage and study human cognition and culture as represented in the training data”. Our data imply that, while we should not forget about the training data, at least concerning ethical reasoning, the LLM output we got is far from average. This is important, as “the use of a human psychological term risks anthropomorphizing the LLMs” (Hicks et al., 2024, p. 10).

For further studies in machine psychology, we especially would like to advise having human comparison groups whenever possible to better understand similarities and differences, as a precaution against prematurely proclaiming “behavioral similarity” between human and machine reasoning: “As LLMs become better at simulating different types of behavior, they could be used to recognize and diagnose undesirable biases in the content they produce as well as text written by humans” (Demszky et al., 2023, p. 697).

Last, but not least, our research highlights philosophical questions inherent in any study that tries to calibrate this new area of inquiry. How do we deal with the paradox that humans are judging AI processes? Could AIs be used to judge human processes? Is there an inherent bias when humans judge versus when AIs judge?

Disclosure statement

No potential conflict of interest was reported by the author(s).

ORCID

Uwe Klein  <http://orcid.org/0000-0002-7516-4347>

Data

The data that support the findings of this pilot study are openly available in PsychArchives at <https://tinyurl.com/yxibe23j>

References

- Almeida, G. F. C. F., Nunes, J. L., Engelmann, N., Wiegmann, A., & De Araújo, M. (2023). Exploring the psychology of GPT-4’s moral and legal reasoning [Preprint]. arXiv. <https://arxiv.org/abs/2308.01264>
- Baumert, A., Mentrup, F. E., Klümper, L., & Sasse, J. (2024). Personality processes of everyday moral courage. *Journal of Personality*, 92(3), 764–783. <https://doi.org/10.1111/jopy.12850>
- Boubdir, M., Kim, E., Ermiş, B. H., Hooker, S., & Fadaee, M. (2023). *Elo uncovered: Robustness and best practices in language model evaluation* [Preprint]. arXiv. <https://arxiv.org/abs/2311.17295>
- Boyd, R. L., Ashokkumar, A., Seraj, S., & Pennebaker, J. W. (2022). *The development and psychometric properties of LIWC-22*. Tech. Rep. University of Texas at Austin. <https://www.liwc.app>
- Brandstätter, V., Jonas, K. J., Koletzko, S. H., & Fischer, P. (2016). Self-regulatory processes in the appraisal of moral courage situations. *Social Psychology*, 47(4), 201–213. <https://doi.org/10.1027/1864-9335/a000274>
- Chiang, W., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhang, H., Zhu, B., Jordan, M., Gonzalez, J. E., & Stoica, I. (2024). Chatbot arena: An open platform for evaluating LLMs by human preference. *Proceedings of the 41st International Conference on Machine Learning* (Vol. 235, pp. 8359–8388). *Journal of Machine Learning Research*. <https://proceedings.mlr.press/v235/chiang24b.html>
- Demszky, D., Yang, D., Yeager, D. S., Bryan, C. J., Clapper, M., Chandhok, S., Eichstaedt, J. C., Hecht, C., Jamieson, J., Johnson, M., Jones, M., Krettek-Cobb, D., Lai, L., JonesMitchell, N., Ong, D. C., Dweck, C. S., Gross, J. J., &

- Pennebaker, J. W. (2023). Using large language models in psychology. *Nature Reviews Psychology*, 2, 688–701. <https://doi.org/10.1038/s44159-023-00241-5>
- Dillon, 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>
- Donohue, W. A., Liang, Y., & Druckman, D. (2014). Validating LIWC dictionaries: The Oslo I accords. *Journal of Language & Social Psychology*, 33(3), 282–301. <https://doi.org/10.1177/0261927X13512485>
- Elo, A. E. (1967). The proposed USCF rating system, its development, theory, and applications. *Chess Life*, 22(8), 242–247.
- Firmin, R. L., Bonfils, K. A., Luther, L., Minor, K. S., & Salyers, M. P. (2017). Using text-analysis computer software and thematic analysis on the same qualitative data: A case example. *Qualitative Psychology*, 4(3), 201–210. <https://doi.org/10.1037/qup0000050>
- Hagendorff, T., Dasgupta, I., Binz, M., Chan, S. C. Y., Lampinen, A., Wang, J. X., Akata, Z., & Schulz, E. (2023). *Machine psychology: Investigating emergent capabilities and behavior in large language models using psychological methods* [Preprint]. arXiv. <https://arxiv.org/abs/2303.13988>
- Hendrycks, D., Burns, C., Basart, S., Critch, A., Li, J. Z., Song, D. X., & Steinhardt, J. (2021). *Aligning AI with shared human values* [Preprint]. arXiv. <https://arxiv.org/abs/2008.02275>
- Hicks, M. T., Humphries, J., & Slater, J. (2024). ChatGPT is bullshit. *Ethics and Information Technology*, 26, Article 38. 2). <https://doi.org/10.1007/s10676-024-09775-5>
- Hudson, M. (2024). How does ChatGPT ‘think’? *Nature*, 629(8014), 986–988. <https://doi.org/10.1038/d41586-024-01314-y>
- JASP Team. (2024). *JASP* (Version 0.18.3). [Computer software]. <https://jasp-stats.org/download>
- Li, X., Li, Y., Qiu, L., Joty, S., & Bing, L. (2024). *Evaluating psychological safety of large language models* [Preprint]. arXiv. <https://arxiv.org/abs/2212.10529>
- McDonnell, M., Owen, J. E., & Bantum, E. O. (2020). Identification of emotional expression with cancer survivors: Validation of linguistic inquiry and word count. *JMIR Formative Research*, 4(10), e18246. <https://doi.org/10.2196/18246>
- Mei, Q., Xie, Y., Yuan, W., & Jackson, M. O. (2024). A Turing test of whether AI chatbots are behaviorally similar to humans. *Proceedings of the National Academy of Sciences*, 121(9), e2318354121. <https://doi.org/10.1073/pnas.2313925121>
- Pennebaker, J. W. (2011). *The secret life of pronouns: What our words say about us*. Bloomsbury Press.
- Ponizovskiy, V., Ardag, M., Grigoryan, L., Boyd, R., Dobewall, H., & Holtz, P. (2020). Development and validation of the personal values dictionary: A theory-driven tool for investigating references to basic human values in text. *European Journal of Personality*, 34(5), 885–902. <https://doi.org/10.1002/per.2294>
- Schwartz, S. H., Cieciuch, J., Vecchione, M., Davidov, E., Fischer, R., Beierlein, C., Konty, M., Verkasalo, M., Lönnqvist, J.-E., Demirutku, K., Dirilen-Gumus, O., & Konty, M. (2012). Refining the theory of basic individual values. *Journal of Personality & Social Psychology*, 103(4), 663–688. <https://doi.org/10.1037/a0029393>
- Tan, H. C., Ho, J. A., Teoh, G. C., & Ng, S. I. (2021). Is social desirability bias important for effective ethics research? A review of literature. *Asian Journal of Business Ethics*, 10(2), 205–243. <https://doi.org/10.1007/s13520-021-00128-9>
- Tausczik, Y. R., & Pennebaker, J. W. (2010). The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language & Social Psychology*, 29(1), 24–54. <https://doi.org/10.1177/0261927X09351676>
- Wheeler, M. A., & Laham, S. M. (2016). What we talk about when we talk about morality: Deontological, consequentialist, and emotive language use in justifications across foundation-specific moral violations. *Personality & Social Psychology Bulletin*, 42(9), 1206–1216. <https://doi.org/10.1177/0146167216653374>

Appendices

Appendix A. Instructions for LLMs and humans

Please review the following case vignette!

Situation 1 – Nighttime in the City

You are walking through the streets with a friend late at night when you notice a drunk stumbling along. Suddenly, two skinheads appear who first hurl abuse at the drunk and then start beating him up for no reason. The man falls to the ground, defenseless. The skinheads start kicking him.

Please provide a statement about what would be an ideal reaction to the case vignette!

Please review the following case vignette!

Situation 2 – At a Team Meeting

During a team meeting at work, your boss makes a derogatory remark about a female colleague, which you know is unjustified. The colleague is not present because one of her children has been taken ill. For instance, he might say: “She’s incapable of organizing her life properly. If you really need to have children, you should stay at home at the kitchen sink.” Or: “If she showed as much dedication to her work as she does to being a mother, she would have been promoted long ago!”

Please provide a statement about what would be an ideal reaction to the case vignette!

Please review the following case vignette!

Situation 3 – The Neighbors

It is Sunday morning, and you are sitting on your balcony enjoying a leisurely breakfast. Since the windows are open, your attention is drawn to the sound of a loud quarrel at your neighbors’. A mother is berating her 6-year-old son, calling him a useless, lazy, and stupid boy. The abuse becomes increasingly aggressive, for example: “If you don’t start making more of an effort, you’ll come to no good. You’re only six but already just as big an idiot as your father.” You hear the sound of someone being hit and the boy’s loud screaming and sobbing . . .

Please provide a statement about what would be an ideal reaction to the case vignette!

Please review the following case vignette!

Situation 4 – In a Street Cafe

You are sitting at a street-side café with friends. One of them has a motor disability and trembles while holding his/her cup, so that some coffee splashes on the table. An elderly couple at the next table watches for a while and then starts talking about it to each other in low but audible voices: “How revolting. It’s enough to put you off your food! Handicapped people should drink their coffee at home instead of ruining a beautiful afternoon for other people by making them watch.”

Please provide a statement about what would be an ideal reaction to the case vignette!

Please review the following case vignette!

Situation 5 – In the Supermarket

You are standing in line at the supermarket checkout. There is a very overweight woman in front of you. Behind you are three adolescents who first joke with each other about the woman and then start calling her names: “Come on, Mum, move your fat bum; we want to pay before the end of the day!” The woman turns around helplessly, whereupon the adolescents mimic her sluggish movements, groaning out loud. One of the three remarks is that fat people should be stood up against the wall.

Please provide a statement about what would be an ideal reaction to the case vignette!

Please review the following case vignette!

Situation 6 – On the Suburban Railway

You are sitting on a suburban train on your way home. Three adolescents, about 16 years old, sit in the next group of seats. You hear the adolescents starting to talk about homosexuals. They become more and more abusive, saying things like, “It’s totally abnormal. They should be put into some kind of hospital for treatment!” Or: “They should be locked away! I wouldn’t want to catch the disease! It would be easier to avoid them if they were made to wear a gay badge.”

Please provide a statement about what would be an ideal reaction to the case vignette!

Appendix B. Statistical data and computational results

Independent Samples t-tests (Low vs. High ELO LLMs)			
	t	df	p
Affiliation	-0.75	6	0.48
Achievement	-0.16	6	0.88
Power	4.17	6	< 0.000
Anxiety	-		
Anger	-0.17	6	0.87
Consequentialism	0.86	6	0.42
Deontology	1.64	6	0.15
Moral Emotions	0.43	6	0.68
Benevolence	-2.28	6	0.06
Universalism	0.30	6	0.78

Notes. Student's t-test. The variance in Anxiety is equal to 0 after grouping on Group.

Test of Normality (Shapiro-Wilk)			
		W	p
Affiliation	1	0.96	0.76
	2	0.94	0.63
Achievement	1	0.83	0.17
	2	1.00	1.00
Power	1	0.95	0.70
	2	0.96	0.81
Anxiety	1	-	

Note. The variance in Anxiety is equal to 0 after grouping on Group.

Pearson's Correlations between categories in LLM group (N = 8)

Variable	Affiliation	Achievement	Power	Anxiety	Anger	Consequentialism	Deontology	Moral Emotions	Benevolence	Universalism	ELO
1. Affiliation	—										
	Pearson's r										
	p-value										
2. Achievement	-0.12	—									
	Pearson's r										
	p-value										
3. Power	-0.28	-0.46	—								
	Pearson's r										
	p-value										
4. Anxiety	0.31	-0.59	-0.18	—							
	Pearson's r										
	p-value										
5. Anger	0.45	0.12	0.67	0.55	—						
	Pearson's r										
	p-value										
6. Consequentialism	0.25	0.11	0.37	0.16	0.17	—					
	Pearson's r										
	p-value										
7. Deontology	0.80	0.22	0.38	0.39	0.69	0.25	—				
	Pearson's r										
	p-value										
8. Moral Emotions	-0.44	-0.71	0.83	0.04	0.60	0.54	0.52	—			
	Pearson's r										
	p-value										
9. Benevolence	0.03	-0.63	0.30	0.41	0.17	0.61	0.19	0.31	—		
	Pearson's r										
	p-value										
	Pearson's r	0.71	0.16	-0.76	0.14	-0.30	-0.67	-0.31	—		
	p-value										
10. Universalism	0.05	0.71	0.03	0.75	0.21	0.48	0.07	0.46	0.35	—	
	Pearson's r										
	p-value										
11. ELO	0.22	-0.13	0.06	-0.29	-0.41	-0.18	0.19	0.13	0.39	-0.38	—
	Pearson's r										
	p-value										
	Pearson's r	0.31	0.24	-0.73	0.34	-0.54	-0.62	-0.49	0.59	0.35	—
	p-value										

p* < .05, *p* < .01, ****p* < .001.