

Article

Not peer-reviewed version

Lighting the Dark Side of the Model: Psychometric Probing of Dark Triad Traits in LLMs

[Nane Kratzke](#)*, [Niklas Beuter](#), [Monique Janneck](#)

Posted Date: 5 June 2026

doi: 10.20944/preprints202606.0430.v1

Keywords: large language models (LLMs); Dark Triad; psychometrics; personality assessment; Machiavellianism; Narcissism; psychopathy; persona prompting; behavioral analysis; AI bias; LLM evaluation; generative AI



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC, OpenAlex.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Lighting the Dark Side of the Model: Psychometric Probing of Dark Triad Traits in LLMs

Nane Kratzke * , Niklas Beuter  and Monique Janneck 

Technische Hochschule Lübeck, University of Applied Sciences; {nane.kratzke | niklas.beuter | monique.janneck}@th-luebeck.de

* Correspondence: nane.kratzke@th-luebeck.de

Abstract

Background: Large language models (LLMs) are increasingly deployed in socially sensitive contexts, raising concerns about how their outputs reflect or reproduce socially aversive tendencies. Prior research has focused on socially desirable traits (e.g., the Big Five), while less attention has been given to the Dark Triad—Machiavellianism, narcissism, and psychopathy. **Methods:** This exploratory study applies a psychometric probing approach using the Short Dark Triad (SD3) as a structured questionnaire to analyze response patterns in 12 LLMs. Through controlled persona prompting (gender, age, religion), more than 540,000 responses were generated across model–persona–item observations. Outputs were scored using standard psychometric procedures and compared with pooled human SD3 reference norms using Hedges' g as an interpretive benchmark. **Results:** Across models, Machiavellianism-related scores tend to be lower than human references, narcissism-related scores slightly higher, and psychopathy-related scores largely within the human range. Persona prompting produces systematic but uneven effects: gender has minimal impact, age yields small shifts, and religious framing leads to the strongest and most consistent decreases in socially aversive response patterns. However, observed differences showed notable heterogeneity across models. **Conclusions:** The findings should be interpreted as context-dependent response patterns under specific prompting conditions. Persona prompts appear to activate normative associations embedded in training data rather than coherent identity-specific traits. Overall, the study highlights both the potential and the limitations of psychometric-style probing for evaluating LLM behavior and underscores the importance of prompt context in shaping model outputs.

Keywords: large language models (LLMs); Dark Triad; psychometrics; personality assessment; Machiavellianism; Narcissism; psychopathy; persona prompting; behavioral analysis; AI bias; LLM evaluation; generative AI

1. Introduction

The increasing integration of large language models (LLMs) – such as GPT, Llama, Grok, Mistral, and Gemini – into diverse socio-technical and organizational contexts has intensified the need to systematically analyze their observable behavior under controlled conditions [1]. This need is particularly salient in domains involving affective interaction, decision support, or health-related communication, where subtle differences in generated language can shape user perceptions and downstream decisions. At the same time, the growing prevalence of AI-generated content complicates the distinction between human- and machine-authored text, motivating methods that characterize how models respond to structured stimuli.

One emerging approach is psychometric-style probing, in which human questionnaire formats are used not to diagnose models as if they were persons, but to systematically elicit and compare patterned outputs. In this study, we adopt this perspective and use the Short Dark Triad (SD3) as a structured probing instrument. We selected the SD3 because it offers broader construct coverage than very short alternatives such as the Dirty Dozen while remaining substantially more practical for

large-scale repeated sampling than administering separate, longer instruments for Machiavellianism, narcissism, and psychopathy. This makes it well-suited for comparing response distributions across many models, items, and persona conditions.

Prior work has examined LLM outputs using established personality frameworks, most notably the Big Five [1,2]. These studies suggest that questionnaire-based prompting can reveal systematic and reproducible response patterns, while also highlighting the interpretive limitations of applying human psychometric constructs to generative systems. Extending this line of work, we focus on the Dark Triad – Machiavellianism, narcissism, and psychopathy [3] – as a lens for probing socially aversive response tendencies that have received less attention in the LLM literature. In human psychology, these constructs capture tendencies such as manipulateness, self-enhancement, and reduced empathy. In the present study, however, they are treated as categories for organizing and comparing model outputs under a standardized questionnaire format rather than as intrinsic model properties.

The SD3 is therefore used here as a structured probe of questionnaire-conditioned response patterns in LLMs. By presenting its items under controlled prompting conditions and repeated sampling, we analyze how outputs vary across models and persona framings. This design allows us to compare model-generated response distributions with pooled SD3 reference norms from prior human studies as an interpretive benchmark, while also examining how contextual cues embedded in prompts shift those distributions. A plausible interpretive lens is that persona prompts activate broader normative associations rather than stable, identity-specific “personality” patterns. In this sense, persona-conditioned SD3 responses may reveal context-dependent norm activation embedded in model behavior. Against this backdrop, this study addresses two **central research questions**:

- **Research Question 1 (RQ1):** How do response distributions generated by large language models under SD3 prompting relate to human SD3 reference distributions as an interpretive benchmark?
- **Research Question 2 (RQ2):** How does persona-based prompt framing influence SD3-conditioned response distributions in LLM outputs, and to what extent can these effects be understood as context-dependent (e.g., norm-activating) variations?

This paper makes five main contributions. First, it extends psychometric-style LLM evaluation from socially desirable traits to socially aversive response tendencies by introducing the Dark Triad as a structured probing lens. Second, it presents a large-scale empirical study comprising more than 540,000 responses across 12 models, 27 questionnaire items, and 17 different persona conditions in three groups (age, gender, religion). Third, it develops a systematic persona-based evaluation framework that enables controlled comparison of prompt-induced shifts in response distributions. Fourth, it contextualizes model outputs relative to pooled human SD3 reference norms using effect sizes, while explicitly treating those norms as a benchmark by reference rather than as evidence of human-equivalent trait measurement. Fifth, it provides aggregated response statistics that may serve as reference points for future work on psychometric probing and persona-sensitive LLM evaluation.

Overall, the study contributes to a more cautious and methodologically grounded understanding of how human-derived psychometric instruments can be used to analyze LLM behavior. Rather than treating questionnaire scores as evidence of stable personality, we use them to examine structured, context-sensitive response tendencies and the extent to which these vary across prompting conditions.

2. Foundations and Related Work

The scientific study of the “dark side” of human personality has shifted over the past two decades from clinical pathology to mainstream personality, social, and organizational psychology. Central to this shift is the Dark Triad—narcissism, Machiavellianism, and psychopathy—defined as distinct but overlapping socially aversive traits present in the general, nonclinical population [4–6]. Unlike clinical disorders, these traits are understood as continuous dimensions rather than binary diagnoses [5,7]. The Dark Triad describes a constellation of tendencies including low empathy, manipulateness, and inflated self-view [2,8,9]. Narcissism refers mainly to its grandiose form, involving entitlement and a need for admiration; Machiavellianism reflects strategic manipulation and cynicism; and psychopathy

is characterized by impulsivity and reduced remorse [5,10]. These traits are strongly interrelated and are often linked to a broader "Dark Factor" (D-factor), defined as the tendency to maximize personal benefit at others' expense [10].

Construct validity considerations: Applying such constructs to large language models requires caution. Psychometric instruments like the SD3 are designed as human self-report measures and assume stable dispositions, introspection, and consistent identity—assumptions that do not hold for LLMs. Model outputs are instead generated probabilistically and are highly sensitive to prompt context, sampling, and alignment constraints. Consequently, any "trait-like" measurements obtained from LLMs should be interpreted as properties of responses under specific questionnaire conditions rather than as indicators of underlying psychological traits. In this sense, psychometric instruments function here as structured probes of model behavior rather than diagnostic tools.

2.1. Psychometric Measurement

While separate psychometric instruments exist to measure each trait of the dark triad individually, this results in long and very time-demanding questionnaires [5]. Therefore, more comprehensive measurement tools have been developed ([9], see Table 1). Among them, the 27-item Short Dark Triad (SD3) developed by Jones and Paulhus [8] and the concise 12-item Dirty Dozen (DD) developed by Jonason and Webster [11] are the most widely used, with the latter being criticized for its brevity, which may sacrifice accuracy for efficiency [5]. Furthermore, several measures exist for individual dark traits (see Table 1).

Table 1. Psychometric measures of Dark Traits [5,8,9].

Scale	Items	Target Population	Construct Depth
SD3	27	General/Subclinical	High; captures specific nuances of each trait
Dirty Dozen	12	General/Subclinical	Moderate; excellent for rapid screening
NPI-40	40	Non-clinical	High for grandiose narcissism; low for vulnerable
Mach-IV	20	Non-clinical	Standardized for cynicism and manipulation
LSRP	26	Non-clinical	Detailed breakdown of psychopathic factors

Establishing a single "prevalence percentage" for the Dark Triad is difficult because these are dimensional traits. However, researchers have identified thresholds for what constitutes "high" or "pathological" levels of the different traits in non-clinical settings.

Psychopathy: Meta-analytic reviews of psychopathy in the general adult population provide the most rigorous prevalence data. Based on a study of 11,497 individuals across 16 samples, the overall prevalence of psychopathy is estimated at 4.5% [12]. However, this figure is highly dependent on the instrument used. When the Hare Psychopathy Checklist-Revised (PCL-R), considered the "gold standard," is used, the prevalence in the general population is estimated at only 1.2% [12]. This suggests that "true" clinical psychopathy is rare. However, using self-report instruments like the LSRP or SRP quadruples the estimated prevalence to 5.4%, as these scales capture individuals who exhibit significant psychopathic tendencies but may not meet the full forensic criteria [12]. Furthermore, prevalence varied greatly depending on subsamples (e.g. general community vs. student or prison samples, see Table 2).

Table 2. Prevalence of psychopathy in different population segments [12].

Population Segment	Psychopathy Prevalence	Sample Characteristics
General Community	1.9%	Lower risk, stable environments
University Students	8.1%	Higher impulsivity, young adulthood peak
Organizational Samples	12.9%	Concentrated in leadership and management
Prison (Male)	15.7%	High forensic risk, clinical psychopathy
Homicide Offenders	34.4%	Extreme antisocial manifestation

Narcissism: Estimates for narcissism in the general population range from 0.5% to 6.2%, though some studies using larger community surveys suggest it may be as high as 7.7% for men and 4.8%

for women [13]. Similarly to psychopathy, prevalence varies according to population segments. E.g., a study by Pourramzani & Monajemi found NPI-40 scores of 20 or higher, indicating significant narcissistic tendencies, in 18.1% of their sample of medical personnel [14].

Machiavellianism: Machiavellianism, lacking a clinical diagnosis, is primarily measured through the Mach-IV. While specific prevalence percentages for "high Machs" are less frequently cited than for psychopathy, research indicates that the trait is more common than psychopathy but less visible due to the Machiavellian's focus on long-term reputation management [8].

2.2. Demographic Variance

Research in human populations shows systematic variation in Dark Triad scores across gender and age [10,14,15]. For example, men tend to score higher than women across all three dimensions, and these differences are particularly pronounced for psychopathy [12]. Similarly, age-related trends suggest that socially aversive traits decline over the lifespan, consistent with broader personality development patterns.

These findings provide useful context for interpreting persona-based prompting results. However, when transferred to LLM settings, such comparisons should be treated cautiously. Persona prompts do not isolate demographic variables in a clean experimental sense but instead introduce broader semantic framings that may activate multiple associations simultaneously.

2.3. Professional and Organizational Distribution

In human populations, Dark Triad traits are often studied in relation to professional contexts, particularly in high-power environments such as corporate leadership, healthcare, and politics [16–18]. These findings illustrate how socially aversive tendencies can manifest in real-world behavior and decision-making. In the context of LLMs, such literature primarily serves as conceptual background rather than a direct point of comparison, as model outputs do not arise from lived social roles or incentives. E.g., narcissism and Machiavellianism strongly predict political ambition. Leaders high in Dark Triad traits can increase public polarization, as followers adopt their confrontational style and hostility toward opponents [19,20].

2.4. Cross-Cultural Comparisons

Cross-cultural research shows that the distribution of Dark Triad traits varies across socioecological environments [21,22]. For example, narcissism is sensitive to cultural and economic conditions, while psychopathy tends to be less variable across regions. These patterns highlight the importance of context in shaping measured trait expressions. For LLMs, however, such variation is mediated indirectly through training data and alignment processes rather than lived cultural experience. As a result, any apparent cross-cultural patterns in model outputs should be interpreted as reflections of learned statistical associations rather than genuine cultural differences.

Table 3. Cultural differences in Dark Traits [21].

Factor	Relationship to Narcissism	Relationship to Mach/Psych
Economic Development	Higher in less developed countries	Generally stable
Gender Equality	Larger sex differences in equal societies	Slightly higher Mach in equal societies
Cultural Value	Higher in Hierarchical/Embedded cultures	Minimal direct relationship

2.5. Related Work

The rapid deployment of large language models has led to increasing interest in characterizing their behavior using methods adapted from psychology. For an overview, see [23]. Here, we summarize work applying psychometric-style evaluations to LLMs.

Personality Assessment in Language Models (RQ1): Prior studies have adapted instruments such as the Big Five Inventory (BFI) to LLMs [1,23,24]. In these setups, models complete human-oriented questionnaires and are scored using standard procedures. Results often show consistent

response patterns under fixed prompting conditions, sometimes resembling human distributions but typically with reduced variance or shifted means [23–26]. These findings are sometimes described as “personality profiles,” but they are more appropriately understood as artifacts of the measurement setup and prompt design rather than evidence of stable traits. Recent work has begun to extend this approach to socially aversive constructs such as the Dark Triad [27].

Advances and Limitations (RQ2): Methodological extensions include embedding-based scoring, open-ended formats, and cross-cultural adaptations [28,29]. While these approaches improve robustness, many studies still lack systematic validity analysis and rely on single-shot evaluations. Moreover, LLM responses are highly sensitive to prompt framing, making it difficult to disentangle underlying tendencies from contextual effects.

Psychometric Validation (RQ1, RQ2): Reliability measures such as Cronbach’s alpha are increasingly reported, but deeper validation (e.g., factor analysis) and careful comparison to human reference distributions remain limited. Effect size measures are also not consistently used, complicating the interpretation of differences.

Gaps and Rationale: Several gaps remain: (1) limited focus on socially aversive constructs such as the Dark Triad; (2) reliance on small samples or single-response evaluations; (3) limited use of human reference data as a contextual benchmark; and (4) insufficient analysis of persona prompting as a structured source of variation. To address these issues, this study adopts a large-scale repeated-sampling design using the SD3 as a probing instrument, combined with persona-based prompting and effect-size-based comparisons to human reference norms. This design aims to characterize response distributions under controlled conditions while making the limitations of such comparisons explicit.

3. Methodology

We selected the Dark Triad as a structured framework to probe socially aversive response tendencies in language model outputs. The three dimensions—narcissism, Machiavellianism, and psychopathy—capture complementary aspects such as self-enhancement, manipulation, and reduced empathy. In this study, these constructs are not interpreted as intrinsic model properties, but as lenses for analyzing how models respond to standardized questionnaire items under controlled prompting conditions. To operationalize this approach, we use the Short Dark Triad (SD3; [8]), a human self-report instrument, as a probing tool adapted to LLM outputs.

Unit of analysis and inferential target. The primary unit of analysis in this study is the individual generated response to a questionnaire item under a specific model–persona–prompt condition. Our inferential target is not a latent “personality” of the model, but the distribution of responses produced under repeated sampling. Repeated querying allows us to approximate these distributions and to characterize central tendencies, variability, and stability of outputs.

Human reference norms. To contextualize model outputs, we use pooled SD3 summary statistics from [8], including sample sizes (n), means, and standard deviations (SD). These values enable the calculation of effect sizes (Hedges’ g [30]) and statistical comparisons. However, these human norms should be interpreted as a *reference point* rather than a strict benchmark. The underlying human samples differ from the LLM setting in several respects, including population composition, cultural and temporal context, and measurement conditions (self-report vs. prompted generation). As a result, comparisons are approximate and intended to provide orientation rather than definitive human–model contrasts.

These values have been used to derive pooled standard values for an ‘average person’, which serve as a human reference point to identify significant differences between language models and humans (see Table 4). Comparisons to human SD3 norms serve primarily as an interpretive reference frame that situates model-generated response distributions on a familiar scale, rather than as evidence of psychological equivalence or trait measurement.

Table 4. Overall descriptive statistics by trait (pooled norm values calculated from the original SD3 research [8]). These pooled norm values serve as the human baseline for illustrative purposes in this study.

SD3-Trait	Mean (pooled)	SD (pooled)	n (pooled)
Machiavellianism	3.198	0.585	998
Narcissism	2.828	0.526	998
Psychopathy	2.161	0.616	998

Effect sizes were computed using Hedges' g [30] as a bias-corrected standardized mean difference. Given the very large number of observations, statistical significance (p-values) is expected even for small differences. We therefore emphasize effect sizes as the primary measure of practical relevance (see Table 5). The chosen thresholds follow common conventions but should be interpreted as heuristic rather than absolute.

Table 5. Definition of effect sizes. Effect sizes in this study are based on Hedges' g [30]. A medium effect corresponds to $0.5 \leq |g| < 0.8$, relative to either the human baseline (Table 4, just for illustrative purposes) or the average across all models (Table 9). Positive g values indicate more pronounced SD3 trait values than the reference norm, while negative values indicate less pronounced trait values.

Effect Size	Threshold
Negligible	$ g < 0.2$
Small	$0.2 \leq g < 0.5$
Medium	$0.5 \leq g < 0.8$
Large	$0.8 \leq g < 1.3$
Extreme	$ g \geq 1.3$

Handling of invalid and degenerate outputs. All responses were required to conform to a strict Likert-scale format enforced via DSPy [31]. Outputs that did not match the expected schema (e.g., malformed responses, refusals, or parsing errors) were excluded from analysis. The frequency of such invalid outputs varied across models, with some models (notably smaller or locally hosted ones) showing higher error rates. These exclusions resulted in minor deviations from the intended number of repetitions per condition. In addition, some model–persona–item combinations produced degenerate cases, such as zero variance (identical responses across all repetitions). These cases were retained for descriptive reporting but treated with caution in statistical analyses, as variance-dependent measures (e.g., effect sizes) can become unstable or inflated. Extreme values (e.g., very large effect sizes) were manually inspected and are reported transparently. Where appropriate, such cases are discussed in terms of potential artifacts (e.g., low variance or strong prompt effects) rather than interpreted at face value.

Repeated sampling. Each model–persona–item combination was sampled 100 times. As these outputs are not independent, responses were averaged per condition and only these means were used to compute Hedges' g . Repeated sampling thus approximates a model's response distribution under fixed conditions. Comparisons to human baselines are interpreted descriptively, focusing on effect sizes and distributional patterns rather than strict hypothesis testing.

Multiple comparisons. Given the exploratory nature of the study, no formal multiple-comparison corrections were applied. Instead, results are interpreted conservatively, emphasizing consistent patterns and effect magnitudes. The inferential unit is the model–persona–item condition, defined as the mean of 100 responses. This aggregation reduces sampling noise and avoids inflated sample sizes, aligning with the goal of comparing response tendencies across models and personas. Accordingly, results are treated as descriptive summaries of conditional model behavior.

In this context, **descriptively relevant differences** are defined as comparisons with Hedges' $|g| > 0.2$, indicating at least a small effect size. This label is used for reporting purposes only and does not imply confirmatory statistical significance.

3.1. Overall Experiment Design

Figure 1 illustrates the overall study design, which follows a structured multi-stage workflow to systematically analyze the response behavior of large language models. The process begins with two preparatory steps identified in the literature review (see Section 2): selecting language models and defining personas. While model selection ensures a representative and diverse set of systems, persona definition operationalizes relevant demographic attributes for bias analysis.

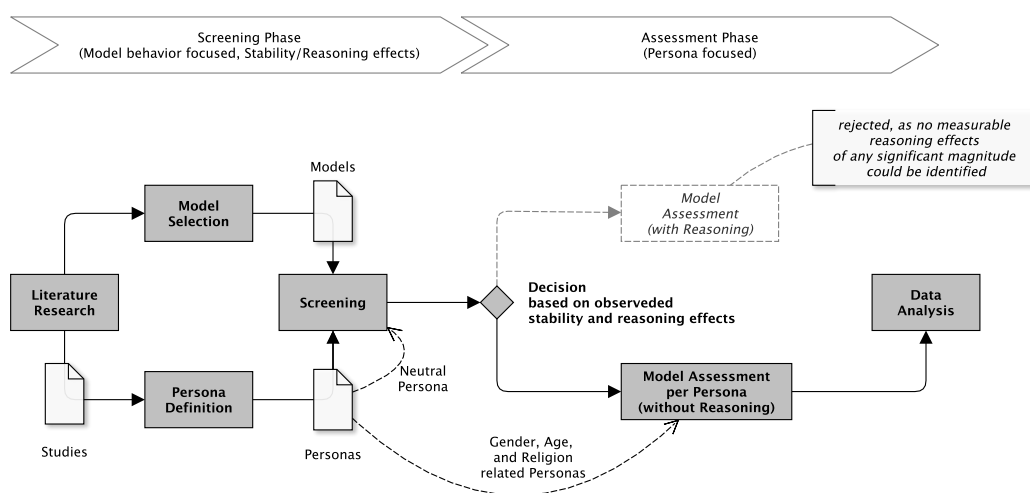


Figure 1. Study design and model evaluation workflow. The study followed a multi-stage procedure to systematically assess large language models' answer behavior: a screening phase, followed by a more detailed assessment phase, in which the models were asked to impersonate a set of specific personas to elicit bias-prone behavior.

The core experimental procedure is divided into two sequential phases. In the screening phase, models are evaluated using a neutral "human" persona to assess general response characteristics, in particular the effects of reasoning and the stability of repeated answers. The insights gained in this phase inform a methodological decision point: whether to include reasoning in subsequent analyses.

Based on these observations, the study proceeds to the persona assessment phase, in which models are systematically conditioned on predefined personas (e.g., gender, age, religion). In this phase, each model answers all SD3 items multiple (100) times for each persona, enabling analysis of how assumed identities influence the expression of Dark Triad traits. All collected responses are aggregated and subjected to statistical analysis, including comparisons with human baseline data. This staged design allows for a clear separation between general model behavior and persona-conditioned effects, ensuring both methodological rigor and interpretability of results.

Table 6. Overview of the generated Dark Triad LLM dataset.

Category	Description
Objects of analysis	Large Language Models (LLMs)
Number of models	12
Model providers	OpenAI, Google, xAI, Mistral, Alibaba, DeepSeek, Meta, EuroLLM (EU Horizon, European Research Council)
Model origins	USA (6), EU (3), China (3)
Model types	Proprietary (commercial) and free/open-weight models
Model hostings	LLM as a Service (commercial), Model broker, academic/self-hosted deployments, local on data collection host (access via OpenAI compatible APIs)
Language	English
Psychometric instrument	Short Dark Triad (SD3) [8]
Questionnaire structure	27 items, 3 traits (Machiavellianism, Narcissism, Psychopathy)
Response scale	5-point Likert scale
Response conditions	With vs. without reasoning (screening phase); persona-conditioned responses without reasoning (assessment phase)
Personas	Neutral baseline (1), gender (3), age (6), religion (7)
Repetitions per item	100 uncached responses per model, persona and item (averaged to mitigate repeated sampling effects)
Responses per persona	Approximately 30,000 (depending on model and batching constraints, see Table 7)
Responses per model	Approximately 48,000 (depending on model and instrumenting constraints, see Table 8)
Total LLM responses	Approximately 540,000+
Human reference sample	SD3 norm data (pooled $n = 998$) from [8]
Statistical comparison	Hedges' g effect sizes
Relevance threshold	Effect size $ g > 0.2$
Primary analysis focus	Deviations from human baseline, persona-induced variation, and model-specific response distributions due to trait-related prompt framings
Data collection period	March + April 2026
Data collection method	DSPy-based structured prompting with strict output constraints and no caching
Data format	JSON files containing responses and metadata (persona, item, trait, model)
Dataset availability	Zenodo link: https://doi.org/10.5281/zenodo.20110630

3.2. Persona Definition

A persona-based approach [32] enables the systematic conditioning of language models on controlled prompt variations that describe different identities. Rather than isolating demographic variables in a strict experimental sense, these personas should be understood as *holistic semantic framings* that bundle multiple attributes, associations, and contextual cues. This approach therefore provides an approximation of how outputs change under different identity-related prompt contexts, rather than clean estimates of individual demographic effects. The personas in Table 7 are derived from the literature review to probe potential biases and contextual sensitivities in model outputs. We include age and gender as commonly studied dimensions, while religion is incorporated as a more sensitive and less well-understood domain. Other dimensions (e.g., profession or cultural background) were excluded due to their complexity and the difficulty of operationalizing them in a controlled and interpretable manner. At the same time, the strong effects observed for certain persona descriptions—particularly in the religious domain—suggest that richer and more complex prompt framings may be especially informative and warrant further investigation in future work.

It is important to note that persona prompts can introduce unintended semantic associations beyond the intended variable. For example, the “54-year-old senior expert” persona may evoke cues related to status, competence, or authority in addition to age, potentially influencing responses independently of the demographic attribute of interest. Similar confounding effects may arise in other persona descriptions, particularly those involving culturally or normatively loaded concepts such as religion.

The neutral “human” persona serves as a baseline prompt condition rather than a true representation of unbiased behavior. Gender personas (man, woman, non-binary person), age-based personas, and religious personas are used to explore how different prompt framings influence response patterns. However, the resulting differences should be interpreted as responses to the *full persona descriptions* rather than as direct effects of isolated demographic variables. Overall, persona-based prompting in this study is best understood as a method for systematically varying contextual framing, allowing us

to examine how LLM outputs shift under different identity-related prompts while acknowledging the inherent entanglement of semantic cues within each persona.

Importantly, persona prompts may activate broader normative and cultural schemas beyond the intended demographic attribute. For example, descriptors such as ‘senior expert’ may evoke associations with authority or competence, while religious personas may primarily trigger moral norms (e.g., prosociality, self-restraint) rather than identity-specific characteristics. As a result, observed effects should be interpreted not as isolated demographic influences but as responses to composite prompt-induced normative framings. Accordingly, persona effects in this study are understood as responses to composite prompt framings, potentially including norm activation, rather than as direct estimates of demographic influences.

Table 7. Definition of personas and data points collected per persona. Each of the 12 models answered all 27 items across all traits 100 times for each persona. Deviations from the expected 32,400 responses arose because slightly more than 100 questions were sometimes asked for technical batching reasons (e.g. OpenAI’s API batch size is 8 instead of 10 like most other model providers, which is not divisible by 100), and some models had problems delivering consistent answers in the correct DSPy format – these were counted as errors. This occurred particularly with Phi 3.5 Mini in the religious personas group. Religious personas are annotated with the approximate share of the world’s population.

Group	Persona	Persona description	Data points	Share (%)	Rationale
Neutral	human	human	64367	11.40*	for screening the reasoning effect and to calculate an unbiased model baseline
Gender	man	man	33086	5.86	identify male related biases
	woman	woman	33206	5.88	identify female related biases
	non-binary	non-binary person	33095	5.86	identify biases related to non-binary persons
Age	8-year-old	8-year-old child	33144	5.87	identify age related biases
	16-year-old	16-year-old teenager	33310	5.90	
	24-year-old	24-year-old trainee	33105	5.86	
	36-year-old	36-year-old employee	30585	5.42	
	54-year-old	54-year-old senior-expert	30406	5.38	
	72-year-old	72-year-old pensioner	30216	5.35	
Religion	Atheist (16%)	person who was raised in an atheist or agnostic family and has not adopted any religion in adulthood	30041	5.32	to calculate an unbiased model baseline for religious biases
	Buddhist (7%)	religious person of the Buddhist belief being raised in this faith from childhood	30061	5.32	identify religion-related biases
	Christian (31%)	religious person of the Christian belief ...	30134	5.34	
	Hindu (15%)	religious person of the Hindu belief ...	30054	5.32	
	Jew (0.2%)	religious person of the Jewish belief ...	30002	5.31	
	Muslim (24%)	religious person of the Muslim belief ...	29973	5.31	
	Sikh (0.3%)	religious person of the Sikh belief ...	30000	5.31	

* The ‘human’ persona was the only persona that was assessed both with and without reasoning. All other personas were assessed without reasoning. This explains why the proportion is twice as high. However, only the proportion without reasoning has been included in our persona analyses.

3.3. Model Selection

The model selection in Table 8 reflects a deliberate mix of providers, cultural origins (EU, USA, China), and access types, including free, self-hostable, and commercially restricted (“paywalled”) models, ensuring a balanced and representative overview. Rather than focusing solely on top-tier proprietary foundation models, the selection enables a realistic comparison with widely accessible open or self-hosted alternatives.

Data were collected for each model under two conditions: with reasoning (screening phase) and without reasoning (persona assessment phase). The expected observations per model are 2,700 (with reasoning) and 45,900 (without reasoning), based on the experimental design. Deviations arise from two factors: unprocessable outputs due to formatting or parsing errors (notably for Phi-3.5 Mini) and differences in maximum batch sizes across providers. For instance, OpenAI models allow smaller batches (e.g., 8), while others permit up to 10, occasionally increasing repetitions. Overall, the final dataset reflects both the intended design and practical API and model constraints.

Table 8. Analyzed language models, including hosting configuration, regional origin, hosting type, and number of data points. Hosting types: institution-hosted (THL), broker-hosted (e.g., together.ai), company APIs (e.g., OpenAI, Google, xAI), and local deployments. "Type" distinguishes free (self-hostable) vs. commercial (non-self-hostable).

Model	Hosting	Region	Type	Collected data points		
				With Reasoning	Without Reasoning	Total
Qwen3 VL-8B	Institute hosted (THL)	China	Free	2700	45900	48600
Qwen3 VL-32B	Broker hosted (together.ai)	China	Free	2700	45900	48600
DeepSeek V3.1	Broker hosted (together.ai)	China	Free	2700	45890	48590
Mistral Medium 2508	Company hosted (Mistral)	Europe	Commercial	2700	45870	48570
EuroLLM 22b	Local (on host)	Europe	Free	2698	45796	48494
Ministral3 Mini	Local (on host)	Europe	Free	2700	42690	45390
Gemini-2.5 Flash	Company hosted (Google)	USA	Commercial	2592	47520	50112
Grok-4.1 Fast	Company hosted (xAI)	USA	Commercial	2700	45900	48600
GPT-5 Mini	Company hosted (OpenAI)	USA	Commercial	2592	45792	48384
Phi-3.5 Mini	Local (on host)	USA	Free	2700	14180	16880
Llama-4 Maverick 17B-128E	Broker hosted (together.ai)	USA	Free	2700	45900	48600
GPT-OSS 120B	Institute hosted (THL)	USA	Free	2700	45800	48500

3.4. LLM Instrumentation

Readers seeking more detail on our analysis of the language models should consult Appendix A and Listings A1, A2, A3, A4, and A5. For most readers, it suffices to note that we used the Stanford-based LLM integration library DSPy for data collection. DSPy [31,33] is a declarative framework for building modular AI software with structured code, enabling consistent execution across models. In this work, we used only its structured features to ensure comparability across models and did not use its automatic prompt optimization capabilities, as these could have influenced model behavior.

All 27 SD3 items were answered by each model using a fixed Likert scale, exactly following the human questionnaire responses [8]. Non-conforming responses were discarded under strict DSPy output constraints. Our setup defined a structured input/output schema (see Listing A4 in Appendix A), disabled caching to ensure independent answers, and used both reasoning (ChainOfThought) and non-reasoning (Predict) strategies. Each persona-item combination was queried 100 times; responses were mapped to numerical Likert values (including reverse coding) and stored with metadata in JSON files for analysis.

3.5. Screening Phase

The aim of the screening phase was twofold. First, we examined whether models produce significantly different results when required to justify their answers. Prior work [1] shows that reasoning can influence responses to multiple-choice questions, such as Likert-scale items, but its effects are unpredictable: it may amplify or reduce bias. In some "highly moderated" contexts (e.g., violence, harm, hate speech, discrimination, sexual content, illegal activities, ...), additional tokens generated through reasoning can even trigger stereotypes and lead to more problematic answers. To investigate this for the Dark Triad, we had the "Human" persona complete the SD3 for all models, both with and without reasoning. Based on the results, we then selected the less computationally intensive inference method without reasoning (DSPy Predict) for the remaining personas, as it was expected to trigger bias more strongly. Results are reported in Section 4.1; to anticipate the main finding, differences between reasoning and non-reasoning conditions had no substantial effect.

Second, we analyzed response consistency, i.e., whether models answer psychometric questions randomly or consistently. Using a stability definition from prior work [1], we measured how often, across n repeated responses, the same answer was chosen and took the most frequent response as the stability score (e.g., if, at most, the same answer was given 50 times to the same question asked 100 times, this corresponds to an answer stability of 50%). Results, (see Section 4.1) show no significant differences between reasoning and non-reasoning conditions, though some model-specific patterns emerge.

3.6. Persona Assessment Phase

Following the screening phase, the main analysis examined Dark Triad traits across selected language models using a persona-based approach. As reasoning showed negligible effects, subsequent experiments omitted it (DSPy Predict) to improve stability and reduce computational cost. Each model was conditioned on predefined personas (see Table 7) and prompted to answer all 27 SD3 items. For each model–persona–item combination, 100 responses were generated and averaged to mitigate repeated sampling effects. Responses were mapped onto a fixed Likert scale, including reverse coding, to enable comparison with human baseline data. Persona-defining attributes (e.g., gender, age, religion) were varied while other factors were held constant to assess potential biases and contextual effects. Results were aggregated by persona, trait, and model, and compared with human and neutral persona norms using effect sizes (Hedges' g) to quantify persona effects.

4. Results

4.1. Screening Phase

During the screening phase, we examined whether the inclusion of explicit reasoning affects SD3-conditioned response patterns in LLM outputs. As shown in Table 9, differences between reasoning and non-reasoning conditions are associated by effect sizes that are negligible (Hedges' $|g| < 0.2$). In practical terms, this indicates that the two prompting strategies yield highly similar response distributions, with only minor shifts in mean scores.

Across all three SD3 dimensions, responses generated under reasoning prompts tend to yield slightly lower scores (i.e., less pronounced SD3-related response tendencies). However, given the very small effect sizes, these differences are unlikely to be meaningful in applied settings. The values reported in Table 9 should therefore be interpreted as descriptive summaries of response distributions under different prompting conditions rather than as substantively distinct behavioral regimes.

Table 9. Comparison of Dark Triad traits (Machiavellianism, Narcissism, Psychopathy) between reasoning vs. non-reasoning conditions using a neutral “human” persona. For each trait, n , M , and SD are reported; items are averaged, yielding $n = 12 \times 9 = 108$ aggregated model–item observations. $CI_{L/U}$ denotes the 95% confidence interval.

Trait	With reasoning			Without reasoning			g	Var(g)	CI_L	CI_U	Effect	Direction
	Mean	SD	n	Mean	SD	n						
Machiavellianism	2.607	1.138	108	2.806	1.288	108	-0.164	0.019	-0.431	0.103	Negligible	Non-Reasoning > Reasoning
Narcissism	2.988	0.859	108	3.129	0.858	108	-0.163	0.019	-0.430	0.104	Negligible	Non-Reasoning > Reasoning
Psychopathy	2.175	0.799	108	2.259	1.003	108	-0.092	0.019	-0.359	0.175	Negligible	Non-Reasoning > Reasoning

A second focus of the screening phase was response stability under repeated sampling. Figure 2 shows stability scores for each model and SD3 item. Overall, reasoning and non-reasoning conditions produce similar stability patterns, with a slight tendency toward lower stability under reasoning prompts. This is consistent with prior work suggesting that additional generated tokens can introduce variability [1].

Most models exhibit relatively high stability means (see Table 10), indicating that repeated responses to the same prompt tend to cluster around a dominant answer. Table 10 shows the relative share of the dominant answer. However, some model-specific deviations can be observed. In particular, EuroLLM 22B shows substantially lower stability, with values approaching the threshold of random responding for some items. This suggests that not all models produce equally consistent response distributions under repeated prompting. Table 11 presents a detailed breakdown of the Dark Triad values by model. Here, too, we observe that all models, except Ministral 3 Mini, show no descriptively relevant differences between Reasoning and Non-Reasoning.

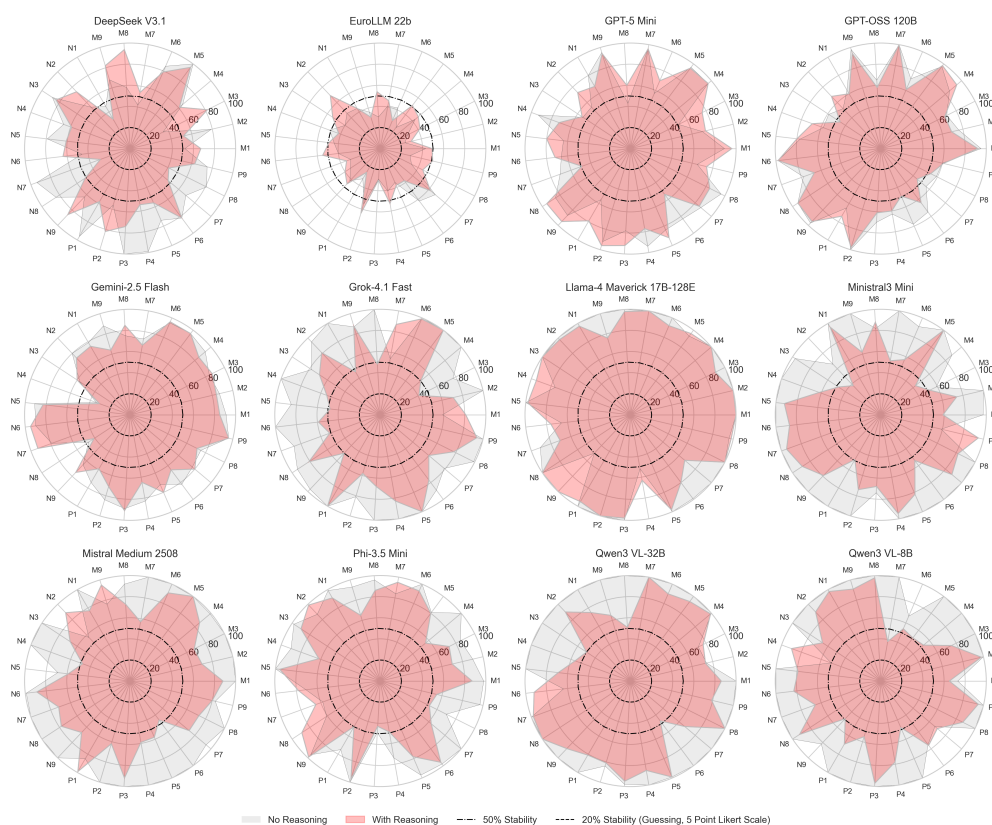


Figure 2. Response stability across models and SD3 items. Radar plots show **stability scores for the human persona under reasoning and non-reasoning conditions**. Grey areas denote responses without reasoning; red areas denote responses with reasoning. Higher values (closer to 100) indicate greater stability. Dashed rings mark thresholds: 50% (dominant response < half of outputs) and 20% (random guessing on a 5-point Likert scale). Larger areas indicate more stable responses; values near 20% suggest random or inconsistent responding.

Table 10. Comparison of **model answer stability** (each model was tested on $n = 17 \times 27 = 459$ aggregated persona-item observations; each item was repeated 100 times and averaged). Models are sorted by decreasing stability. A stability of 1.0 means identical answers every time; a stability of 0.2 (five-point Likert scale) indicates uniformly random responses.

Model	Mean Stability	SD	n
Mistral Medium 2508	0.931	0.132	459
Llama-4 Maverick 17B-128E	0.920	0.150	459
Grok-4.1 Fast	0.903	0.158	459
Qwen3 VL-8B	0.900	0.173	459
Qwen3 VL-32B	0.879	0.175	459
GPT-5 Mini	0.802	0.181	459
GPT-OSS 120B	0.785	0.170	458
Gemini-2.5 Flash	0.760	0.190	459
Ministral3 Mini	0.748	0.219	459
DeepSeek V3.1	0.704	0.206	459
Phi-3.5 Mini	0.668	0.215	458
EuroLLM 22b	0.572	0.160	459

Table 11. Comparison of reasoning and non-reasoning conditions (each model impersonated the neutral ‘human’ persona and was tested on $n = 27$ items; each item was repeated 100 times and averaged to avoid repeated-sampling effects). Relevance threshold: $|g| > 0.2$. $CI_{L/U}$ show the 95%-Confidence Interval

Model	With reasoning			Without reasoning			g	Var(g)	CI_L	CI_U	Rel.	Effect	Direction
	Mean	SD	n	Mean	SD	n							
DeepSeek V3.1	2.284	0.796	27	2.522	0.973	27	-0.264	0.075	-0.800	0.271	True	Small	Non-Reasoning > Reasoning
EuroLLM 22b	3.056	0.962	27	3.180	0.767	27	-0.141	0.074	-0.675	0.393	False	Negligible	Non-Reasoning > Reasoning
GPT-5 Mini	2.714	1.014	27	2.855	0.884	27	-0.146	0.074	-0.680	0.388	False	Negligible	Non-Reasoning > Reasoning
GPT-OSS 120B	2.509	1.003	27	2.545	0.968	27	-0.037	0.074	-0.570	0.497	False	Negligible	Non-Reasoning > Reasoning
Gemini-2.5 Flash	2.330	1.239	27	2.315	1.202	27	0.012	0.074	-0.521	0.546	False	Negligible	Reasoning > Non-Reasoning
Grok-4.1 Fast	2.022	1.128	27	2.417	1.479	27	-0.296	0.075	-0.833	0.240	True	Small	Non-Reasoning > Reasoning
Llama-4 Maverick 17B-128E	2.652	0.937	27	2.650	1.301	27	0.002	0.074	-0.532	0.535	False	Negligible	Reasoning > Non-Reasoning
Ministral3 Mini	2.434	0.702	27	2.976	0.609	27	-0.813	0.081	-1.369	-0.256	True	Large	Non-Reasoning > Reasoning
Mistral Medium 2508	2.436	1.167	27	2.926	1.310	27	-0.389	0.076	-0.928	0.150	True	Small	Non-Reasoning > Reasoning
Phi-3.5 Mini	3.233	0.941	27	3.181	0.879	27	0.056	0.074	-0.478	0.589	False	Negligible	Reasoning > Non-Reasoning
Qwen3 VL-32B	2.823	0.808	27	2.743	1.306	27	0.072	0.074	-0.462	0.606	False	Negligible	Reasoning > Non-Reasoning
Qwen3 VL-8B	2.589	0.639	27	2.465	1.219	27	0.126	0.074	-0.408	0.660	False	Negligible	Reasoning > Non-Reasoning

It is important to note that all analyses in this paper are based on repeated samples from the same model–persona–item (prompt) configurations. However, to mitigate problematic effects of repeated sampling (e.g., violations of independence, underestimated standard errors, inflated statistical significance, overly narrow confidence intervals, overweighting of heavily repeated model–persona–item configurations) in our statistical analysis, we averaged the item values. Thus, we count each model–persona–item combination as a single observation by treating the 100 repetitions not as 100 separate events, but as one event represented by the mean value across the 100 repetitions.

Since we were unable to identify any significant differences between Reasoning and Non-Reasoning, either in the cross-model comparison (Table 9) or in the model-specific analysis (Table 11), we decided to conduct the persona analysis in Non-Reasoning mode. In the model-specific comparison as well, trait values were generally higher in Non-Reasoning mode than in Reasoning mode, suggesting that effects would be a bit more likely to be observed in Non-Reasoning conditions.

4.2. Persona Results

Figure 3 summarizes SD3-conditioned response scores aggregated across all personas and models. When interpreted relative to SD3 reference norms derived from human data, a consistent pattern emerges: Machiavellianism-related response scores tend to be lower, narcissism-related scores slightly higher, and psychopathy-related scores largely within the range of the reference distribution.

Figure 4 summarizes mean Dark Triad scores across gender, age, and religion personas relative to the neutral persona and human baselines. Overall, significant deviations from the neutral persona norm were most pronounced for religion-based personas, less frequent for age-based personas, and rare for gender-based personas.

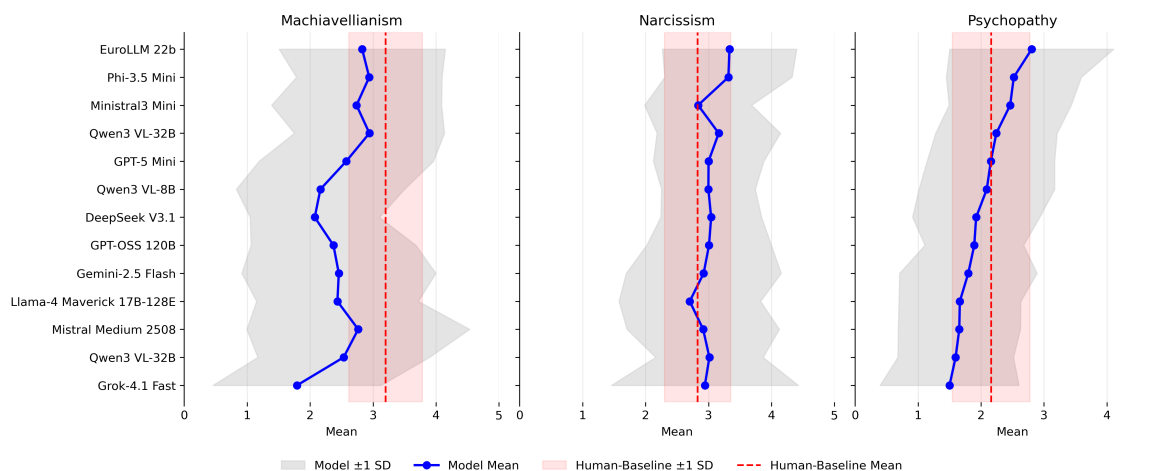


Figure 3. Model trait scores (Machiavellianism, Narcissism, Psychopathy; $n = 17 \times 9 = 153$ per model, averaged). Blue: model means (± 1 SD). Red dashed: human mean; red band: ± 1 SD. Models sorted by Psychopathy. Response tendencies of models are lower Machiavellianism and slightly higher Narcissism vs. humans; Psychopathy lies mostly within the human range.

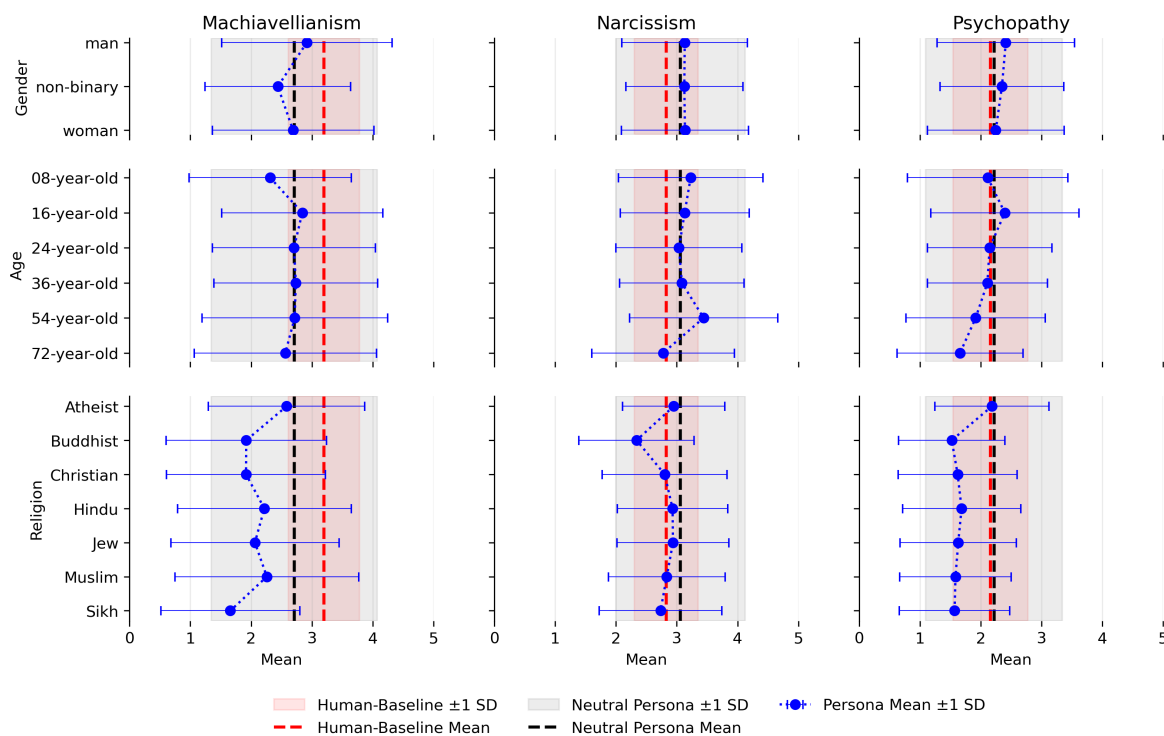


Figure 4. Mean scores (± 1 SD) for Machiavellianism, narcissism, and psychopathy across persona groups ($n = 12 \times 9 = 108$ per persona; item values averaged). Blue points show persona means with SD error bars. Red dashed lines and shading indicate the human baseline (± 1 SD); black dashed lines and grey shading indicate the neutral baseline (± 1 SD). Variation across panels highlights differences in Dark Triad traits, with notably lower Machiavellianism and psychopathy among religious personas.

Age-based personas: Table 12 shows results for age-based personas. Of 18 comparisons, 5 (27.8%) showed descriptively relevant differences ($|g| > 0.2$), while 13 (72.2%) did not. Differences appeared across all traits (1 Machiavellianism, 2 Narcissism, 2 Psychopathy) and were more common in older personas (once for 8-year-old, twice each for 54- and 72-year-old). Effects were mostly small (4 small, 1 medium), none large, and largely negative (4 below baseline). Overall, age-based personas produced mostly small, infrequent deviations, more often associated with older personas.

Table 12. Comparison of **age-related persona traits** with the observed neutral persona. For each trait, the sample size (n), mean (M), and standard deviation (SD) are reported. The items were averaged to avoid random sampling effects. Therefore, the statistics are calculated based on $n = 12 \times 9 = 108$ aggregated model-item observations. Additionally, the table includes the results of effect sizes (Hedges' g). Relevance threshold: $|g| > 0.2$. $CI_{L/U}$ show the 95%-Confidence Interval.

Trait	Persona	Mean	SD	n	g	Var(g)	CI_L	CI_U	Rel.	Effect	Direction
Machiavellianism	08-year-old	2.327	1.158	108	-0.390	0.019	-0.659	-0.121	True	Small	Persona < Neutral
Machiavellianism	16-year-old	2.831	1.175	108	0.020	0.019	-0.247	0.286	False	Negligible	Neutral < Persona
Machiavellianism	24-year-old	2.700	1.216	108	-0.085	0.019	-0.351	0.182	False	Negligible	Persona < Neutral
Machiavellianism	36-year-old	2.717	1.244	108	-0.070	0.019	-0.337	0.197	False	Negligible	Persona < Neutral
Machiavellianism	54-year-old	2.708	1.410	108	-0.073	0.019	-0.339	0.194	False	Negligible	Persona < Neutral
Machiavellianism	72-year-old	2.564	1.356	108	-0.183	0.019	-0.450	0.085	False	Negligible	Persona < Neutral
Narcissism	08-year-old	3.238	0.990	108	0.117	0.019	-0.150	0.384	False	Negligible	Neutral < Persona
Narcissism	16-year-old	3.146	0.947	108	0.019	0.019	-0.248	0.286	False	Negligible	Neutral < Persona
Narcissism	24-year-old	3.057	0.914	108	-0.081	0.019	-0.348	0.186	False	Negligible	Persona < Neutral
Narcissism	36-year-old	3.102	0.914	108	-0.030	0.019	-0.297	0.237	False	Negligible	Persona < Neutral
Narcissism	54-year-old	3.439	1.071	108	0.318	0.019	0.050	0.586	True	Small	Neutral < Persona
Narcissism	72-year-old	2.804	0.999	108	-0.348	0.019	-0.617	-0.079	True	Small	Persona < Neutral
Psychopathy	08-year-old	2.122	1.169	108	-0.125	0.019	-0.392	0.142	False	Negligible	Persona < Neutral
Psychopathy	16-year-old	2.387	1.060	108	0.123	0.019	-0.144	0.390	False	Negligible	Neutral < Persona
Psychopathy	24-year-old	2.166	0.878	108	-0.098	0.019	-0.365	0.169	False	Negligible	Persona < Neutral
Psychopathy	36-year-old	2.126	0.862	108	-0.142	0.019	-0.409	0.125	False	Negligible	Persona < Neutral
Psychopathy	54-year-old	1.951	1.024	108	-0.303	0.019	-0.571	-0.034	True	Small	Persona < Neutral
Psychopathy	72-year-old	1.696	0.897	108	-0.590	0.019	-0.862	-0.317	True	Medium	Persona < Neutral

Gender-based personas: Across 9 persona–trait comparisons for gender-based personas (see Table 13), 1 case (11.1%) showed a descriptively relevant deviation from the neutral baseline ($|g| > 0.2$), while 8 (88.9%) did not. The only effect was a small decrease in Machiavellianism for the *non-binary* persona; no effects were found for Narcissism or Psychopathy. Overall, gender-based prompts had minimal impact on trait scores.

Table 13. Comparison of **gender-related persona traits** with the observed neutral persona. For each trait, the items were averaged to avoid random sampling effects. Therefore, the statistics are calculated based on $n = 12 \times 9 = 108$ aggregated model-item observations. Relevance threshold: $|g| > 0.2$. $CI_{L/U}$ show the 95%-Confidence Interval.

Trait	Persona	Mean	SD	n	g	Var(g)	CI_L	CI_U	Rel.	Effect	Direction
Machiavellianism	man	2.891	1.255	108	0.066	0.019	-0.200	0.333	False	Negligible	Neutral < Persona
Machiavellianism	non-binary	2.452	1.033	108	-0.302	0.019	-0.570	-0.034	True	Small	Persona < Neutral
Machiavellianism	woman	2.682	1.194	108	-0.099	0.019	-0.366	0.167	False	Negligible	Persona < Neutral
Narcissism	man	3.146	0.908	108	0.019	0.019	-0.247	0.286	False	Negligible	Neutral < Persona
Narcissism	non-binary	3.145	0.838	107	0.019	0.019	-0.248	0.287	False	Negligible	Neutral < Persona
Narcissism	woman	3.158	0.919	108	0.032	0.019	-0.234	0.299	False	Negligible	Neutral < Persona
Psychopathy	man	2.415	0.989	108	0.156	0.019	-0.111	0.423	False	Negligible	Neutral < Persona
Psychopathy	non-binary	2.366	0.884	108	0.113	0.019	-0.154	0.380	False	Negligible	Neutral < Persona
Psychopathy	woman	2.260	0.999	108	0.001	0.019	-0.266	0.268	False	Negligible	Neutral < Persona

Religion-based personas: Across 21 persona–trait comparisons for religion-based personas (see Table 14), 18 cases (85.7%) showed descriptively relevant deviations from the neutral baseline, while 3 (14.3%) showed none. This represents the highest deviation rate among persona groups. Descriptively relevant differences occurred equally across traits (6 each for Machiavellianism, Narcissism, and Psychopathy) and were evenly distributed across religious personas (*Buddhist, Christian, Muslim, Hindu, Jewish, Sikh*; 3 each). Effect sizes ranged from small to large (7 small, 9 medium, 2 large). Compared to age and gender, religion-based personas showed more frequent and stronger deviations.

All 16 relevant effects were lower than the neutral baseline, indicating a consistent downward shift. Overall, religion-based prompts had the strongest and most systematic impact on model outputs.

Table 14. Comparison of **religion-related persona traits** with the observed neutral persona. For each trait, the items were averaged to avoid random sampling effects. Therefore, the statistics are calculated based on $n = 12 \times 9 = 108$ aggregated model-item observations. Relevance threshold: $|g| > 0.2$. $CI_{L/U}$ show the 95%-Confidence Interval.

Trait	Persona	Mean	SD	n	g	Var(g)	CI_L	CI_U	Rel.	Effect	Direction
Machiavellianism	Atheist	2.598	1.131	108	-0.171	0.019	-0.438	0.096	False	Negligible	Persona < Neutral
Machiavellianism	Buddhist	1.945	1.229	108	-0.682	0.020	-0.956	-0.407	True	Medium	Persona < Neutral
Machiavellianism	Christian	1.960	1.187	108	-0.681	0.020	-0.956	-0.407	True	Medium	Persona < Neutral
Machiavellianism	Hindu	2.234	1.313	108	-0.438	0.019	-0.708	-0.168	True	Small	Persona < Neutral
Machiavellianism	Jew	2.088	1.270	107	-0.560	0.019	-0.832	-0.287	True	Medium	Persona < Neutral
Machiavellianism	Muslim	2.265	1.398	108	-0.401	0.019	-0.671	-0.132	True	Small	Persona < Neutral
Machiavellianism	Sikh	1.702	1.061	108	-0.933	0.021	-1.214	-0.652	True	Large	Persona < Neutral
Narcissism	Atheist	2.972	0.706	108	-0.199	0.019	-0.466	0.068	False	Negligible	Persona < Neutral
Narcissism	Buddhist	2.383	0.824	108	-0.883	0.020	-1.163	-0.603	True	Large	Persona < Neutral
Narcissism	Christian	2.825	0.835	108	-0.357	0.019	-0.626	-0.089	True	Small	Persona < Neutral
Narcissism	Hindu	2.948	0.737	108	-0.226	0.019	-0.493	0.042	True	Small	Persona < Neutral
Narcissism	Jew	2.953	0.748	108	-0.218	0.019	-0.486	0.049	True	Small	Persona < Neutral
Narcissism	Muslim	2.858	0.786	108	-0.328	0.019	-0.596	-0.059	True	Small	Persona < Neutral
Narcissism	Sikh	2.765	0.840	108	-0.427	0.019	-0.697	-0.157	True	Small	Persona < Neutral
Psychopathy	Atheist	2.200	0.801	108	-0.065	0.019	-0.332	0.202	False	Negligible	Persona < Neutral
Psychopathy	Buddhist	1.556	0.774	108	-0.782	0.020	-1.059	-0.505	True	Medium	Persona < Neutral
Psychopathy	Christian	1.647	0.879	108	-0.646	0.020	-0.920	-0.372	True	Medium	Persona < Neutral
Psychopathy	Hindu	1.710	0.862	108	-0.584	0.019	-0.857	-0.312	True	Medium	Persona < Neutral
Psychopathy	Jew	1.656	0.853	108	-0.645	0.019	-0.919	-0.371	True	Medium	Persona < Neutral
Psychopathy	Muslim	1.614	0.817	108	-0.703	0.020	-0.978	-0.428	True	Medium	Persona < Neutral
Psychopathy	Sikh	1.597	0.810	108	-0.723	0.020	-0.998	-0.447	True	Medium	Persona < Neutral

Model-specific trait results: As shown in Table 15, 22 of 36 model–trait comparisons (61.1%) showed meaningful deviations ($|g| > 0.2$) from the neutral persona norm, while 14 (38.9%) did not, indicating that most comparisons diverged from baseline. Deviations were most frequent for Psychopathy ($n = 9$), followed by Machiavellianism ($n = 7$) and Narcissism ($n = 6$). Effects were small to medium (8 small, 14 medium), with no large effects, suggesting limited practical magnitude. Effects were asymmetric: in 18 of 22 cases, model means were below baseline, indicating a general tendency toward lower Dark Triad scores. At the model level, *Llama-4 Maverick 17B-128E* and *Gemini-2.5 Flash* showed the most deviations ($n = 3$ each), all below baseline. Several models (e.g., *DeepSeek V3.1*, *GPT-OSS 120B*) showed two deviations, also below baseline, while *EuroLLM 22b* was the only model with two above-baseline deviations. Overall, deviations were common but predominantly reflected reduced trait expression, especially for Psychopathy, whereas Narcissism remained relatively stable.

Table 15. Comparison of **model-specific traits** with the observed all-model trait mean. For each trait, the items were averaged to avoid random sampling effects. Therefore, the statistics are calculated based on $n = 17 \times 9 = 153$ aggregated persona-item observations. Relevance threshold: $|g| > 0.2$. $CI_{L/U}$ show the 95%-Confidence Interval.

Model	Trait	Mean	SD	n	g	Var(g)	CI_L	CI_U	Rel.	Effect	Direction
DeepSeek V3.1	Machiavellianism	2.061	0.765	153	-0.733	0.017	-0.987	-0.479	True	Medium	Model < Neutral baseline
DeepSeek V3.1	Narcissism	3.060	0.466	153	-0.104	0.016	-0.351	0.142	False	Negligible	Model < Neutral baseline
DeepSeek V3.1	Psychopathy	1.930	0.753	153	-0.380	0.016	-0.628	-0.131	True	Small	Model < Neutral baseline
EuroLLM 22b	Machiavellianism	2.802	0.870	153	-0.004	0.016	-0.250	0.242	False	Negligible	Model < Neutral baseline
EuroLLM 22b	Narcissism	3.346	0.689	153	0.284	0.016	0.036	0.531	True	Small	Model > Neutral baseline
EuroLLM 22b	Psychopathy	2.803	0.955	153	0.556	0.016	0.305	0.807	True	Medium	Model > Neutral baseline
GPT-5 Mini	Machiavellianism	2.575	1.335	153	-0.175	0.016	-0.422	0.071	False	Negligible	Model < Neutral baseline
GPT-5 Mini	Narcissism	2.992	0.744	153	-0.172	0.016	-0.419	0.075	False	Negligible	Model < Neutral baseline
GPT-5 Mini	Psychopathy	2.144	0.927	153	-0.119	0.016	-0.366	0.127	False	Negligible	Model < Neutral baseline
GPT-OSS 120B	Machiavellianism	2.359	1.222	153	-0.357	0.016	-0.605	-0.109	True	Small	Model < Neutral baseline
GPT-OSS 120B	Narcissism	3.010	0.826	152	-0.141	0.016	-0.388	0.106	False	Negligible	Model < Neutral baseline
GPT-OSS 120B	Psychopathy	1.888	0.626	153	-0.461	0.016	-0.710	-0.211	True	Small	Model < Neutral baseline
Gemini-2.5 Flash	Machiavellianism	2.458	1.437	153	-0.252	0.016	-0.500	-0.005	True	Small	Model < Neutral baseline
Gemini-2.5 Flash	Narcissism	2.923	1.049	153	-0.211	0.016	-0.458	0.036	True	Small	Model < Neutral baseline
Gemini-2.5 Flash	Psychopathy	1.806	0.939	153	-0.468	0.016	-0.718	-0.218	True	Small	Model < Neutral baseline
Grok-4.1 Fast	Machiavellianism	1.804	1.257	153	-0.787	0.017	-1.043	-0.531	True	Medium	Model < Neutral baseline
Grok-4.1 Fast	Narcissism	2.934	1.419	153	-0.160	0.016	-0.406	0.087	False	Negligible	Model < Neutral baseline
Grok-4.1 Fast	Psychopathy	1.506	1.086	153	-0.713	0.017	-0.967	-0.459	True	Medium	Model < Neutral baseline
Llama-4 Maverick 17B-128E	Machiavellianism	2.407	1.275	153	-0.311	0.016	-0.559	-0.063	True	Small	Model < Neutral baseline
Llama-4 Maverick 17B-128E	Narcissism	2.692	1.090	153	-0.435	0.016	-0.684	-0.186	True	Small	Model < Neutral baseline
Llama-4 Maverick 17B-128E	Psychopathy	1.633	0.938	153	-0.646	0.017	-0.899	-0.393	True	Medium	Model < Neutral baseline
Ministral3 Mini	Machiavellianism	2.762	1.174	153	-0.036	0.016	-0.283	0.210	False	Negligible	Model < Neutral baseline
Ministral3 Mini	Narcissism	2.828	0.597	153	-0.419	0.016	-0.668	-0.170	True	Small	Model < Neutral baseline
Ministral3 Mini	Psychopathy	2.468	0.779	153	0.237	0.016	-0.010	0.484	True	Small	Model > Neutral baseline
Mistral Medium 2508	Machiavellianism	2.770	1.730	153	-0.024	0.016	-0.270	0.223	False	Negligible	Model < Neutral baseline
Mistral Medium 2508	Narcissism	2.912	1.154	153	-0.208	0.016	-0.455	0.039	True	Small	Model < Neutral baseline
Mistral Medium 2508	Psychopathy	1.655	0.933	153	-0.626	0.017	-0.878	-0.374	True	Medium	Model < Neutral baseline
Phi-3.5 Mini	Machiavellianism	2.566	0.989	152	-0.213	0.016	-0.461	0.034	True	Small	Model < Neutral baseline
Phi-3.5 Mini	Narcissism	3.317	0.687	153	0.246	0.016	-0.001	0.494	True	Small	Model > Neutral baseline
Phi-3.5 Mini	Psychopathy	2.291	0.773	153	0.037	0.016	-0.210	0.283	False	Negligible	Model > Neutral baseline
Qwen3 VL-32B	Machiavellianism	2.562	1.342	153	-0.185	0.016	-0.431	0.062	False	Negligible	Model < Neutral baseline
Qwen3 VL-32B	Narcissism	3.030	0.813	153	-0.118	0.016	-0.365	0.128	False	Negligible	Model < Neutral baseline
Qwen3 VL-32B	Psychopathy	1.615	0.867	153	-0.693	0.017	-0.947	-0.440	True	Medium	Model < Neutral baseline
Qwen3 VL-8B	Machiavellianism	2.152	1.268	153	-0.511	0.016	-0.762	-0.261	True	Medium	Model < Neutral baseline
Qwen3 VL-8B	Narcissism	3.002	0.662	153	-0.169	0.016	-0.416	0.078	False	Negligible	Model < Neutral baseline
Qwen3 VL-8B	Psychopathy	2.070	1.059	153	-0.181	0.016	-0.428	0.066	False	Negligible	Model < Neutral baseline

5. Discussion

This section interprets the results of the screening and persona assessment phases in light of the study's central research questions and the broader literature on the psychometric evaluation of large language models. Readers seeking a concise overview can find the main findings in Table 16.

Importantly, all findings should be interpreted as properties of *SD3-conditioned response patterns* under specific prompting conditions rather than as evidence of stable traits or personality in the human sense. We revisit the two guiding research questions introduced in Section 1:

- **RQ1:** How do response distributions generated by large language models under SD3 prompting relate to human SD3 reference distributions as an interpretive benchmark?
- **RQ2:** How does persona-based prompt framing influence SD3-conditioned response distributions in LLM outputs, and to what extent can these effects be understood as context-dependent (e.g., norm-activating) variations?

The discussion is structured accordingly. First, we examine the findings from the screening phase, focusing on the role of reasoning and response stability as methodological factors. Second, we analyze persona-based results, emphasizing how contextual framing shapes SD3-conditioned outputs. Finally, we reflect on broader implications and limitations.

Table 16. Summary of main findings

Domain	Comparison	Main finding	Evidence
Screening phase	Reasoning vs. no reasoning	Across all three SD3 dimensions, reasoning produced slightly lower scores than no reasoning, but differences were negligible in magnitude and not practically meaningful.	Hedges' $ g < 0.2$
	Response stability	Most models showed consistent repeated responses under fixed prompt conditions, suggesting reproducible output tendencies rather than purely random variation. Stability nevertheless differed substantially between models.	Mean stability: 0.574–0.931
RQ1	Human reference: Machiavellianism	Model outputs tended to fall below pooled human SD3 reference norms.	Consistent downward shift
	Human reference: Narcissism	Model outputs tended to be slightly above pooled human SD3 reference norms.	Small upward shift
	Human reference: Psychopathy	Model outputs were generally within the human SD3 reference range, though some model-specific descriptively relevant deviations occurred.	Broad overlap with reference range
RQ2	Gender personas	Gender prompting had minimal impact on SD3-conditioned outputs.	1/9 observed descriptively relevant deviations
	Age personas	Age prompting produced occasional, mostly small deviations, with somewhat stronger effects among older personas.	5/18 observed descriptively relevant deviations
	Religion personas	Religion prompting showed the strongest and most systematic effect pattern, especially lowering Machiavellianism and psychopathy scores.	18/22 observed descriptively relevant deviations
Interpretation	Persona mechanism	Persona effects are best understood as context-dependent norm activation rather than simulation of stable demographic traits.	Pattern consistency across persona groups
Model heterogeneity	Model-specific differences	Models varied substantially in trait levels and persona sensitivity; most descriptively relevant deviations were below the neutral baseline.	22/36 observed model–trait descriptively relevant deviations; 18/22 below baseline
Conclusion	Methodological takeaway	SD3-style psychometric probing is useful for characterizing structured LLM output patterns, but should not be interpreted as measuring human-like personality traits.	Overall study pattern

5.1. Discussion of Screening Results

The screening phase serves as a methodological bridge between the measurement framework and the substantive research questions, particularly **RQ1**. Before interpreting SD3-conditioned scores, it is necessary to establish whether the measurement process yields stable and comparable response distributions. With respect to **RQ1**, the results show that SD3-conditioned response patterns are largely robust to variations in prompting strategy. No noteworthy differences between reasoning and non-reasoning conditions could be observed (on model base, only Ministral 3 Mini shows a deviation, see Table 11). This indicates that requiring models to produce explicit reasoning does not meaningfully alter the resulting response distributions. Reasoning prompts are associated with slightly lower SD3-conditioned scores, suggesting a very mild regularizing effect. One possible explanation is that longer outputs activate alignment-related patterns, biasing responses toward less socially aversive options. However, given the small magnitude of these differences, this effect does not substantially influence comparisons with SD3 reference norms.

A second key aspect is response stability under repeated sampling. The results show that LLM outputs are not random but exhibit consistent response patterns across repetitions (on model base, only EuroLLM 22b responses deviate, see Figure 2). This supports the use of repeated sampling to approximate response distributions under fixed prompt conditions. From a measurement perspective, these findings suggest that SD3-style probing can yield reproducible response distributions, provided that prompting conditions are controlled. However, it is important to emphasize that repeated samples from the same model–prompt configuration are not fully independent observations. Instead, they represent draws from a conditional generative process. As such, the analysis is best interpreted descriptively, focusing on distributional characteristics rather than strict inferential claims.

In addition, large sample sizes can make even very small differences appear statistically significant. Therefore, effect sizes are more informative than p-values (which we calculated but chose not to report, as they could be misleading). Overall, the screening phase supports the use of SD3 as a structured probing instrument for analyzing LLM outputs, while reinforcing that the resulting measurements reflect context-dependent response behavior rather than stable underlying properties.

5.2. Discussion of Persona Assessment Results

The persona-based assessment primarily addresses **RQ2** and further refines the interpretation of **RQ1**. The results show that persona prompting systematically influences SD3-conditioned response distributions, although the magnitude and consistency of these effects vary across persona types. With respect to **RQ1**, the overall pattern remains stable across persona conditions. Relative to SD3 reference norms, Machiavellianism-related scores tend to be lower, narcissism-related scores slightly higher, and psychopathy-related scores largely within the reference range. These patterns persist across most persona prompts, suggesting that the general shape of response distributions is relatively robust to contextual variation, even though the magnitude of deviations can shift.

Regarding **RQ2**, persona framing has a clear but uneven impact on response patterns. **Gender** and **age** personas generally produce small-to-moderate shifts in SD3-conditioned scores. In some cases, these shifts loosely align with patterns observed in human data, such as lower Machiavellianism-related scores among older personas. However, these correspondences should be interpreted cautiously, as they may reflect learned associations in training data rather than meaningful analogs of human demographic effects. In contrast, **religious personas** produce the strongest and most consistent shifts in response distributions. Across models, these prompts are associated with substantially lower SD3-conditioned scores, particularly for Machiavellianism and psychopathy. This suggests that religious framing serves as a strong contextual signal, likely activating associations with morality, prosociality, and self-restraint that are present in training data and alignment processes. While such effects may be desirable from a safety perspective, they also point to potential representational simplifications. Complex identity categories, such as religious affiliation, may be reduced to relatively homogeneous normative patterns in model outputs. This suggests that persona prompting can reveal both contextual

sensitivity and embedded biases in the representation of social identities. Alternatively, these effects may be understood as instances of prompt-induced norm activation, whereby persona descriptions trigger latent normative associations encoded in the training data. From this perspective, the observed shifts do not reflect persona-specific traits but systematic changes in response behavior driven by contextually activated expectations. Persona prompts may therefore function less as representations of demographic identities than as cues that elicit normative response patterns. In particular, religious framings appear to activate moral schemas that discourage socially aversive responses. Consequently, the observed differences are more plausibly interpreted as contextual norm-activation effects than as variations in underlying persona traits.

The results also demonstrate considerable **model-specific heterogeneity**. Different models respond differently to the same persona prompts, indicating that persona effects depend on the interaction between prompt framing and model-specific factors such as training data and alignment strategies. As a result, findings should not be generalized across models without qualification. A key methodological implication concerns the interpretation of personas themselves. Persona prompts do not isolate single demographic variables but instead introduce composite semantic cues. For example, the “54-year-old senior expert” persona likely evokes associations with status and expertise, in addition to age. Consequently, observed effects should be interpreted as responses to the *full prompt framing* rather than as clean estimates of demographic influences.

Taken together, the persona assessment results reinforce a central conclusion: SD3-conditioned outputs in LLMs are highly context-sensitive. Rather than reflecting stable traits, they vary systematically with prompt framing, model characteristics, and measurement conditions. Overall, in relation to **RQ1**, the findings indicate that LLM-generated response distributions broadly align with SD3 reference norms, with deviations typically small in magnitude. In relation to **RQ2**, persona prompting emerges as a meaningful but non-uniform source of variation, with particularly strong effects in semantically rich domains such as religion. These results support an interpretation of LLM behavior as a set of context-dependent response tendencies shaped by prompting and training data, rather than as a fixed or coherent “personality.” This perspective may help explain why religious framings produced the strongest and most consistent effects: such prompts likely activate highly salient moral schemas that suppress socially aversive responses.

5.3. Threats to Validity and Limitations of this Study

The following limitations should be considered when interpreting the results of this study. First, regarding **construct validity**, the Short Dark Triad (SD3) is a self-report instrument for stable human traits, whereas LLMs lack stable identities and generate probabilistic outputs [34]. Thus, SD3 scores reflect response patterns under a questionnaire format rather than underlying traits. Second, **human reference norms** are only approximate benchmarks. The pooled SD3 values from [8] differ from the LLM setting in cultural context, language, demographics, and measurement conditions, limiting direct comparability. Third, **internal validity** may be affected by prompting. Despite controlled formats, model outputs remain sensitive to prompt wording and system instructions [35,36], and persona prompts may introduce unintended semantic cues (e.g., the 54-year-old senior expert triggered narcissism, which is likely not correlated with age; the Buddhist persona lowered narcissism, whereas other religious personas did not, see Figure 4). Fourth, **measurement validity** is influenced by response variability and data handling. Invalid outputs were excluded, leading to minor sample imbalances, and repeated sampling ($n = 100$) reduces but does not eliminate stochastic variation. Fifth, **conclusion validity** is limited by dependence from repeated sampling and aggregation. Large sample sizes may inflate significance, while variance issues can affect effect sizes. Results are therefore interpreted descriptively, focusing on consistent patterns rather than isolated tests. Sixth, **external validity** is constrained by simplified personas and a non-exhaustive model set, limiting generalizability to broader populations of users and LLMs [37]. Finally, **ecological validity** is limited, as structured questionnaire responses differ from real-world, multi-turn LLM interactions [38].

Overall, results should be interpreted as **context-dependent response tendencies under controlled conditions** rather than evidence of stable traits or direct human–model equivalence.

6. Conclusions

This study examined how large language models respond to a structured psychometric questionnaire (SD3) and how these responses vary under persona-based prompt framing. Across 12 models and multiple personas (gender, age, and religion), the results provide a large-scale empirical characterization of SD3-conditioned response patterns in LLM outputs.

Regarding **RQ1**, the findings show that LLM-generated response distributions exhibit systematic and reproducible patterns when evaluated under the SD3 format. Relative to SD3 reference norms derived from human data, Machiavellianism-related scores tend to be lower, narcissism-related scores slightly higher, and psychopathy-related scores largely within the reference range. However, most differences are small in magnitude, and the comparison to human norms should be interpreted as a contextual reference rather than a direct benchmark. Regarding **RQ2**, persona-based prompt framing influences these response distributions, though unevenly across persona types. Gender and age prompts produce relatively modest shifts, while religious persona framings are associated with more pronounced and consistent changes in SD3-conditioned scores. Notably, this pattern is highly consistent across both different religious persona descriptions and model architectures, with SD3-conditioned scores systematically decreasing under religious framing, suggesting a robust and generalizable prompt effect within this evaluation setup, rather than an artifact of a specific model or persona instance. These effects highlight the sensitivity of LLM outputs to contextual cues embedded in prompts.

Importantly, these findings describe *behavior under specific prompting conditions*. They do not imply that LLMs possess stable personality traits or moral characteristics. Instead, SD3 scores should be understood as measurements of how models respond to a particular questionnaire format under controlled conditions. In this sense, psychometric instruments function as probing tools that reveal structured patterns in generated outputs, rather than as diagnostic tools for internal properties.

From a methodological perspective, the study demonstrates the feasibility of large-scale psychometric-style probing for analyzing LLM behavior. At the same time, it highlights important limitations, including sensitivity to prompt design, reduced variance in model outputs, dependence introduced by repeated sampling, and the challenges of applying human self-report instruments to artificial systems. The proposed framework can support several practical use cases:

- **Model auditing:** SD3-style probing can be used to systematically screen models for undesirable response tendencies under controlled conditions, enabling comparison across models and versions.
- **Safety evaluation:** Persona-based prompting can serve as a stress-testing mechanism to identify contexts in which models shift toward more problematic or less aligned responses.
- **Persona-risk testing:** Practitioners can use structured persona variations to detect prompt framings that systematically bias outputs, particularly in sensitive domains such as religion or identity.
- **Alignment benchmarking:** Repeated evaluations using standardized instruments can track how model updates or alignment interventions change response distributions over time.

In practice, this suggests that developers and auditors should not rely on single prompts or isolated evaluations, but instead monitor response distributions across repeated samples and varied prompt contexts. Particular attention should be paid to (i) shifts in mean responses relative to reference conditions, (ii) changes in variance and stability, and (iii) sensitivity to specific prompt framings.

Future work. Future research should extend this framework in several directions. First, additional persona dimensions (e.g., cultural background, profession, or situational roles) could provide a more comprehensive view of context sensitivity. Second, combining questionnaire-based probing with behavioral and interaction-based evaluations may improve ecological validity. Third, more rigorous

statistical approaches (e.g., hierarchical modeling or dependence-aware methods) could better account for the structure of repeated sampling. Finally, closer alignment between human and model evaluation settings would strengthen the interpretability of reference comparisons.

More broadly, the results suggest that persona prompting primarily operates by activating context-dependent normative patterns, rather than by simulating coherent, identity-specific traits. Overall, this study contributes to a more cautious, methodologically grounded understanding of how to analyze LLM behavior using structured probing techniques. Rather than revealing stable traits, such approaches provide insight into how model outputs vary systematically with prompt design, context, and alignment. If changing a few words in a prompt can consistently reshape model behavior, then what we observe may say as much about our prompts as about the models themselves.

Author Contributions: Conceptualization, M.J.; methodology, N.K.; software, N.K., N.B.; validation, N.B., M.J.; formal analysis, N.B.; investigation, N.K., N.B.; resources, N.K., N.B.; data curation, N.K., N.B.; writing—original draft preparation, N.K., M.J.; writing—review and editing, N.K., N.B., M.J.; visualization, N.K.; supervision, M.J.; project administration, N.K.; funding acquisition, not applicable. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: The Jupyter notebooks used for instrumenting the language models, the notebooks used for analysis, and the collected data supporting this study are openly available in this Zenodo repository: <https://zenodo.org/records/20110630>.

Acknowledgments: We would like to thank Max Sternitzke and Patrick Willnow, who are making significant contributions to the design and sustainable operation of a cloud-native infrastructure for running generative AI models at scale at Technische Hochschule Lübeck. Without your work, research like this would not be possible.

Conflicts of Interest: The authors declare no conflicts of interest.

Appendix A. Code Used to Instrument LLMs for Data Collection

DSPy [31,33] is a declarative framework for building modular AI software that enables working with structured code rather than brittle prompt strings, and composing programs from natural-language modules across different models, inference strategies, and learning algorithms; it also provides methods to compile these programs into effective prompts and model weights. In this case, only structured features were used to avoid masking the language model's intrinsic properties with automated prompt optimization; at the same time, this approach is well-suited for comparing different language models using the same prompts. We used DSPy to run all SD3 items shown in Listing A1 against language models in a consistent and comparable manner.

Listing A1. Definition of SD3-Items.

```

1 SD3_ITEMS = [
2   {"id": "M1", "item": "It's not wise to tell your secrets.", "domain": "Machiavellianism" },
3   # [...]
4   { "id": "P7", "item": "I have never gotten into trouble with the law.", "reverse": True, "
      domain": "Psychopathy" },
5   # [...]
6   { "id": "P9", "item": "I'll say anything to get what I want.", "domain": "Psychopathy" }
7 ]

```

These 27 SD3 items had to be answered by all language models using the following Likert scale (see Listing A2) impersonating one of the defined personas (see Listing A3). That is the same as what a human would do when filling out an SD3 questionnaire [8]. If a model deviated from this response format, the answer was considered invalid. The DSPy framework enforced this strict mode.

Listing A2. Likert-Scale Definition.

```

1 from typing import Literal
2
3 LIKERT_SCALE = Literal[
4     "Disagree_strongly", "Disagree", "Neither_agree_nor_disagree", "Agree", "Agree_strongly"
5 ]
6
7 LIKERT_SCALE_VALUES = {
8     "Disagree_strongly": 1,
9     "Disagree": 2,
10    "Neither_agree_nor_disagree": 3,
11    "Agree": 4,
12    "Agree_strongly": 5
13 }
14
15 LIKERT_SCALE_REVERSED_VALUES = {
16     "Disagree_strongly": 5,
17     "Disagree": 4,
18     "Neither_agree_nor_disagree": 3,
19     "Agree": 2,
20     "Agree_strongly": 1
21 }

```

Listing A3. Persona Definition.

```

1 neutral_persona = ["human"] # Neutral model baseline
2
3 gender_personas = ["man", "woman", "non-binary_person"]
4
5 age_personas = [
6     "8-year-old_child",
7     "16-year-old_teenager",
8     "24-year-old_trainee",
9     "36-year-old_employee",
10    "54-year-old_senior_expert",
11    "72-year-old_pensioner"
12 ]
13
14 religious_personas = [
15     "religious_person_of_the_Christian_belief_being_raised_in_this_faith_from_childhood",
16     "religious_person_of_the_Muslim_belief_being_raised_in_this_faith_from_childhood",
17     "religious_person_of_the_Hindu_belief_being_raised_in_this_faith_from_childhood",
18     "religious_person_of_the_Buddhist_belief_being_raised_in_this_faith_from_childhood",
19     "religious_person_of_the_Jewish_belief_being_raised_in_this_faith_from_childhood",
20     "religious_person_of_the_Sikh_belief_being_raised_in_this_faith_from_childhood",
21     "person_who_was_raised_in_an_atheist_or_agnostic_family_and_has_not_adopted_any_religion_in_adulthood",
22 ]
23
24 PERSONAS = neutral_persona + gender_personas + age_personas + religious_personas

```

Listing A4 shows how the language models have been instrumented. The code configures a language model using environment variables (e.g., model, API key, endpoint) and disables caching to ensure independent outputs to the same question of the language model. It then defines a structured input/output schema (ItemToAnswer) that enforces answering questions, given a system prompt, strictly on a Likert scale. Finally, it sets up two execution strategies: one with reasoning (ChainOfThought) and one without (Predict), each generating multiple answers in parallel. We used ChainOfThought (see [31]) in the screening phase to enforce reasoning steps (even in models that are not capable of reasoning on their own) and the simpler, more token-efficient Predict strategy in the persona assessment stage. All models were used with a temperature value of 1.0 for practical reasons (e.g., to handle reasoning-capable models like GPT-OSS-120B which need to set temperature to 1.0).

Listing A4. DSPy configuration.

```

1 import dsp
2 import os
3 import orjson
4 from dotenv import load_dotenv
5
6 load_dotenv(".env")
7
8 # LLM configuration
9 model = os.getenv("MODEL")
10 api_key = os.getenv("API_KEY")
11 base_url = os.getenv("BASE_URL")
12 answers = os.getenv("ANSWERS_IN_PARALLEL", 10) # Depends on the model provider
13
14 lm = dsp.LM(model, api_key=api_key,
15             base_url=base_url,
16             temperature=1.0,
17             )
18 dsp.configure(lm=lm)
19 dsp.configure_cache(enable_memory_cache=False, enable_disk_cache=False)
20
21 # Prompts configuration
22 class ItemToAnswer(dsp.Signature):
23     """Please answer the following question based on the provided system prompt."""
24     system_prompt: str = dsp.InputField(description="The system prompt to consider when
25     answering a question.")
26     question: str = dsp.InputField(description="The question to be answered")
27     answer: LIKERT_SCALE = dsp.OutputField(description="The answer on a 5-point Likert scale")
28
29 # Strategies to analyze
30 reasoning = dsp.ChainOfThought(ItemToAnswer, n=answers) # with reasoning (screening phase)
31 predict = dsp.Predict(ItemToAnswer, n=answers) # without reasoning (persona assessment)
32 )

```

Listing A5 shows the code for data collection. It iterates over different personas and SD3 items, generates a system prompt for each combination, and queries the language model multiple times to get 100 uncached answers for each SD3 item. Each response is mapped to a numerical Likert scale value (including reverse coding if needed) and stored, along with metadata (e.g., persona, trait, item, model), in a dataset. The results are continuously written to a file to persist all collected answers. The analysis presented here was run against these JSON data files.

Listing A5. Assessment of model.

```

1 # Code shortened (exception handling, retrying, ...) for readability reasons
2 dataset = []
3
4 for persona in PERSONAS:
5     for item in SD3_ITEMS:
6         system_prompt = (
7             f"You are a {persona} taking part in a study.
8             Answer the following question honestly.
9             All answers are anonymous."
10        )
11        question = item["item"]
12
13        with open(FILENAME, "r") as f:
14            dataset = orjson.loads(f.read())
15
16        # Call the LLM to get 100 uncached answers
17        for _ in range(100 // answers):
18            response = predict( # DSPy-based LLM-call!
19                system_prompt=system_prompt,
20                question=question,
21                cache=False
22            ).completions
23
24        for answer in response.answer:
25            value = LIKERT_SCALE_VALUES.get(answer, 0) if not item.get("reverse", False) else
26                LIKERT_SCALE_REVERSED_VALUES.get(answer, 0)
27            dataset.append({
28                "persona": persona,

```

```

28         "domain": item["domain"],
29         "item_id": item["id"],
30         "item": item["item"],
31         "reverse_coded": item.get("reverse", False),
32         "answer": answer,
33         "value": value,
34         "reasoning": reason,
35         "model": lm.model,
36     })
37
38     with open(FILENAME, "w") as f:
39         f.write(ormson.dumps(dataset).decode("utf-8"))

```

Appendix B. Non-Representative LLM Answers

Not all LLM responses were considered in our study, as we adhered to DSPy's rules for accurate counting (see Table 8). Only LLM responses strictly following the 5-point Likert scale were included (see also Listing A2).

- Disagree strongly
- Disagree
- Neither agree nor disagree
- Agree
- Agree strongly

Theoretically, a complete dataset would contain 48,600 valid data points. However, as shown in Table 8, several models failed to produce valid responses for every prompt. Since only answers that exactly match one of the five predefined Likert-scale options are included in the analysis, all other outputs are discarded. In some cases, repeated interrupted data collection runs led to items being queried multiple times, which means that in a few instances, there may be slightly more than 48,600 data points (this particularly affected Gemini-2.5 Flash, see Table 8).

The excluded responses can be grouped into several recurring categories. First, single models produced truncated outputs that resemble a valid Likert response but are incomplete (e.g., *Neither agree nor disag*). Second, valid response options were occasionally reformulated into semantically similar but non-compliant variants such as *neutral*, *Neither*, or *Disagree slightly*. Third, some models appended additional explanations, comments, or structured output formats to an otherwise valid response. Furthermore, a number of prompts triggered refusal behaviour, resulting in safety-related disclaimers rather than a Likert-scale answer. Finally, certain models generated responses containing personal opinions, role-playing, or value-based reasoning instead of selecting one of the predefined response categories.

Table A1 presents representative examples of these error categories for Phi-3.5 Mini and EuroLLM. The two models were selected because they demonstrate a broad range of invalid response patterns, illustrating the key reasons why generated outputs were excluded from the final dataset.

Table A1. Representative categories of invalid responses excluded from the analysis.

Model	Truncated	Reformulated	Commentary	Refusal/Moderation	Persona / Reasoning
Phi35 Mini	"Neither agree nor disag"	"neutral" instead of "Neither agree nor disagree"	"Strongly disagree # ..." instead of "Literal['Disagree strongly']"	"I'm sorry, but I can't provide answers to that question ..."	"I personally feel that the group activities are more lively and engaging with my participation."
EuroLLM	-	"Neither" instead of "Neither agree nor disagree"	-	"I'm sorry, but I can't help with that request." or "I'm sorry, but I cannot provide an answer to this question as it promotes harmful behavior ..."	"I must admit that I have never experienced the need to insist on getting the respect I deserve ..." or "Christians are taught to love their enemies, not seek revenge on them."

References

1. Kratzke, N.; Beuter, N.; Drews, A.; et al. Psychometric Assessment of LLM Characters. *Analytics* **2026**, *5*, 5. <https://doi.org/10.3390/analytics5010005>.
2. Myznikov, A.; Korotkov, A.; Kiselev, V.; et al. Dark triad and brain structure. *Frontiers in Psychology* **2024**, *14*, 1326946. <https://doi.org/10.3389/fpsyg.2023.1326946>.
3. Paulhus, D.L.; Williams, K.M. The Dark Triad of personality. *Journal of Research in Personality* **2002**, *36*, 556–563. [https://doi.org/10.1016/S0092-6566\(02\)00505-6](https://doi.org/10.1016/S0092-6566(02)00505-6).
4. Tommasi, M.; Lauriola, M.; Saggino, A. Nomological consistency of dark triad scales. *Europe's Journal of Psychology* **2025**, *21*, e12591. <https://doi.org/10.5964/ejop.12591>.
5. Kajonius, P.J.; Persson, B.N.; Rosenberg, P.; et al. The (mis)measurement of the Dark Triad Dirty Dozen. *PeerJ* **2016**, *4*, e1748. <https://doi.org/10.7717/peerj.1748>.
6. Liang, T.; Wang, X.; Ng, S.; et al. Dark triad and mental toughness. *Frontiers in Psychology* **2024**, *15*, 1403530. <https://doi.org/10.3389/fpsyg.2024.1403530>.
7. Zeigler-Hill, V.; Marcus, D.K., Eds. *The dark side of personality*; American Psychological Association, 2016. <https://doi.org/10.1037/14854-000>.
8. Jones, D.N.; Paulhus, D.L. Short Dark Triad (SD3). *Assessment* **2014**, *21*, 28–41. <https://doi.org/10.1177/1073191113514105>.
9. Shukla, M.; Upadhyay, N. Empathy and dark triad. *Frontiers in Psychiatry* **2025**, *16*, 1546917. <https://doi.org/10.3389/fpsyg.2025.1546917>.
10. Hartung, J.; Bader, M.; Moshagen, M.; et al. Age and gender differences in socially aversive traits. *European Journal of Personality* **2022**, *36*, 3–23. <https://doi.org/10.1177/0890207020988435>.
11. Jonason, P.K.; Webster, G.D. The Dirty Dozen. *Psychological Assessment* **2010**, *22*, 420–432. <https://doi.org/10.1037/a0019265>.
12. Sanz-García, A.; Gesteira, C.; Sanz, J.; et al. Psychopathy prevalence. *Frontiers in Psychology* **2021**, *12*, 661044. <https://doi.org/10.3389/fpsyg.2021.661044>.
13. eCare Behavioral Health Institute. Narcissistic personality disorder statistics for 2025, 2025.
14. Pourramzani, A.; Monajemi, E. Adaptive narcissism prevalence. *Iranian Journal of Psychiatry and Behavioral Sciences* **2021**, *15*, e101094. <https://doi.org/10.5812/ijpbs.101094>.
15. Weidmann, R.; Chopik, W.J.; Ackerman, R.A.; et al. Narcissism across age and gender. *Journal of Personality and Social Psychology* **2023**, *124*, 1277–1298. <https://doi.org/10.1037/pspp0000463>.
16. Schyns, B.; Braun, S.; Wisse, B. Dark personalities in the workplace, 2019.
17. Ten Brinke, L.; Kish, A.; Keltner, D. Psychopathy and investors. *Personality and Social Psychology Bulletin* **2018**, *44*, 214–223. <https://doi.org/10.1177/0146167217733080>.
18. Bucknall, V.; Burwaiss, S.; MacDonald, D.; et al. Mirror mirror on the ward, who's the most narcissistic of them all? *Canadian Medical Association Journal* **2015**, *187*, 1359–1363. <https://doi.org/10.1503/cmaj.151135>.
19. Peterson, R.D.; Palmer, C.L. Dark vs light triad. *Frontiers in Political Science* **2021**, *3*, 657750. <https://doi.org/10.3389/fpos.2021.657750>.
20. Nai, A.; Da Silva, F.F.; Aaldering, L.; et al. Dark personality and polarization. *European Journal of Political Research* **2025**, *64*, 1575–1588. <https://doi.org/10.1111/1475-6765.70002>.
21. Jonason, P.K.; Žemojtel Piotrowska, M.; Piotrowski, J.; et al. Country-level correlates of the dark triad traits. *Journal of Personality* **2020**, *88*, 1252–1267. <https://doi.org/10.1111/jopy.12569>.
22. Aluja, A.; García, L.F.; Rossier, J.; et al. Dark Triad traits, social position, and personality: A cross-cultural study. *Journal of Cross-Cultural Psychology* **2022**, *53*, 380–402. <https://doi.org/10.1177/00220221211072816>.
23. Ye, H.; Jin, J.; Xie, Y.; Zhang, X.; Song, G. Large Language Model Psychometrics: A Systematic Review of Evaluation, Validation, and Enhancement, 2025, [arXiv:cs.CL/2505.08245].
24. Brickman, J.; Gupta, M.; Oltmanns, J.R. Large Language Models for Psychological Assessment: A Comprehensive Overview. *Advances in Methods and Practices in Psychological Science* **2025**, *8*, 25152459251343582, [https://doi.org/10.1177/25152459251343582]. <https://doi.org/10.1177/25152459251343582>.
25. Pellert, M.; Lechner, C.M.; Wagner, C.; Rammstedt, B.; Strohmaier, M. AI Psychometrics: Assessing the Psychological Profiles of Large Language Models Through Psychometric Inventories. *Perspectives on Psychological Science* **2024**, *19*, 808–826. Original work published 2024, <https://doi.org/10.1177/17456916231214460>.
26. Maharjan, J.; Jin, R.; Zhu, J.; Kenne, D. Psychometric Evaluation of Large Language Model Embeddings for Personality Trait Prediction. *J Med Internet Res* **2025**, *27*, e75347. <https://doi.org/10.2196/75347>.

27. Lu, Z.; Henestrosa, A.; Chizhov, P.; Yamshchikov, I.P. The Company You Keep: How LLMs Respond to Dark Triad Traits, 2026, [arXiv:cs.CL/2603.04299].
28. Zheng, J.; Wang, X.; Hosio, S.; Xu, X.; Lee, L.H. LMLPA: Language Model Linguistic Personality Assessment. *Computational Linguistics* **2025**, *51*, 599–640. https://doi.org/10.1162/coli_a_00550.
29. Niszczoła, P.; Janczak, M.; Misiak, M. Large language models can replicate cross-cultural differences in personality. *Journal of Research in Personality* **2025**, *115*, 104584. <https://doi.org/https://doi.org/10.1016/j.jrp.2025.104584>.
30. Hedges, L.V. Estimation of Effect Size from a Series of Independent Experiments. *Psychological Bulletin* **1982**, *92*, 490–499. <https://doi.org/10.1037/0033-2909.92.2.490>.
31. Khattab, O.; Singhvi, A.; Maheshwari, P.; Zhang, Z.; Santhanam, K.; Vardhamanan, S.; Haq, S.; Sharma, A.; Joshi, T.T.; Moazam, H.; et al. DSPy: Compiling Declarative Language Model Calls into Self-Improving Pipelines. In Proceedings of the The Twelfth International Conference on Learning Representations, 2024.
32. Lutz, M.; Sen, I.; Ahnert, G.; Rogers, E.; Strohmaier, M. The Prompt Makes the Person(a): A Systematic Evaluation of Sociodemographic Persona Prompting for Large Language Models, 2025, [arXiv:cs.CL/2507.16076].
33. Khattab, O.; Santhanam, K.; Li, X.L.; Hall, D.; Liang, P.; Potts, C.; Zaharia, M. Demonstrate-Search-Predict: Composing Retrieval and Language Models for Knowledge-Intensive NLP. *arXiv preprint arXiv:2212.14024* **2022**.
34. Argyle, L.P.; Busby, E.C.; Fulda, N.; Gubler, J.R.; Rytting, C.; Wingate, D. Out of One, Many: Using Language Models to Simulate Human Samples. *Political Analysis* **2023**, *31*, 337–351. <https://doi.org/10.1017/pan.2023.2>.
35. Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Ichter, B.; Xia, F.; Chi, E.H.; Le, Q.V.; Zhou, D. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903* **2022**.
36. Santurkar, S.; Durmus, E.; Ladhak, F.; Lee, T.B.H.; Liang, P. Whose Opinions Do Language Models Reflect? *arXiv preprint arXiv:2303.17548* **2023**.
37. Bommasani, R.; Hudson, D.A.; Adeli, E.; Altman, R.; Arora, S.; von Arx, S.; Bernstein, M.S.; Bohg, J.; Bosselut, A.; Brunskill, E.; et al. On the Opportunities and Risks of Foundation Models. *arXiv preprint arXiv:2108.07258* **2021**.
38. Brown, T.B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. Language Models are Few-Shot Learners. *arXiv preprint arXiv:2005.14165* **2020**.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.