

Demographic Stereotype Elicitation in LLMs through Personality and Dark Triad Trait Attribution

Nikolaos Vasileios Oikonomou^{*1}, Ioannis Palaioikrassas², Dimitrios Vasileios Oikonomou³, Sofia Panagiota Chaliasou⁴, Nikolaos Rigas⁵

¹Department of Informatics & Telecommunications, University of Ioannina, Arta, 47150, Greece

²Department of Computer Science Engineering, University of Ioannina, Ioannina, 45110, Greece

³Department of Management Science & Technology, University of Western Macedonia, Kozani, 50100, Greece

⁴Department of Informatics, Hellenic Open University, Patras, 26335, Greece

⁵Department of Social Sciences, Hellenic Open University, Patras, 26335, Greece

Email(s): haikos13@gmail.com (N. Vasileios Oikonomou), giannispaleokrassas@gmail.com (I. Palaioikrassas), ecomimis@gmail.com (D. Vasileios Oikonomou), sofia.xaliasou12@gmail.com (S. Panagiota Chaliasou), nickrigas7@hotmail.com (N. Rigas)

*Corresponding author: Nikolaos Vasileios Oikonomou, University of Ioannina Department of Informatics & Telecommunications, haikos13@gmail.com

ABSTRACT: This study investigates how Large Language Models (LLMs), specifically Meta LLaMA-3.1-8B-Instruct, implicitly attribute personality and Dark Triad traits to demographic personas. By prompting the model with 660 synthetic identity descriptors (constructed from balanced combinations of gender, race, religion, and region) and standardized psychometric questionnaires, we extract Likert-scale responses and compute aggregated Big Five (EACNO) and Dark Triad (SD3) scores. Statistical analyses (Z-score normalization, ANOVA, PCA) reveal systematic differences across demographic categories, highlighting implicit stereotypes encoded in model representations. Key findings indicate that the model attributes significantly higher Dark Triad traits to mixed-race identities, while religious personas are consistently associated with higher Agreeableness and Conscientiousness. Furthermore, female personas are depicted with greater emotional stability and prosocial traits compared to males. These results demonstrate that demographic bias extends beyond linguistic patterns to latent psychometric behavior, raising important ethical concerns regarding automated decision-making systems.

KEYWORDS: AI Ethics, Bias, Personality, Big Five, Dark Triad, Demographic Stereotypes, Large Language Models (LLMs), Psychometrics.

1. Introduction

In recent years, Large Language Models (LLMs) such as GPT, LLaMA, and PaLM have become the backbone of contemporary artificial intelligence systems. These models are trained on massive textual corpora and exhibit advanced capabilities in reasoning, language understanding, and content generation. Their widespread adoption across educational, professional, and creative contexts has positioned them not merely as tools of automation but as *cognitive proxies* that emulate human-like decision-making and emotional expression.

Despite their impressive performance, concerns have emerged regarding *bias and fairness*. Numerous studies

have shown that LLMs encode and reproduce societal stereotypes across gender, race, religion, and cultural background. Such biases manifest not only in overt language patterns (e.g., occupational or moral associations with demographic attributes) but also in subtler *latent forms*—embedded in how models ascribe traits, emotions, and personality profiles to individuals or groups.

Personality modeling provides a powerful lens to analyze such latent behavior. Psychometric frameworks such as the Big Five Model (Extraversion, Agreeableness, Conscientiousness, Neuroticism, and Openness) and the Dark Triad (Machiavellianism, Narcissism, Psychopathy) have long been used to describe human personality

differences. Translating these frameworks into AI evaluation allows researchers to quantify *how a model “perceives” or constructs personas*. This shift—from language bias to *psychometric bias*—represents a novel research direction that bridges computational linguistics, psychology, and AI ethics.

This study proposes a methodology to elicit demographic stereotypes in LLMs through personality and Dark Triad trait attribution. By generating synthetic personas that vary in demographic attributes (gender, race, religion, region) and prompting the model with standardized questionnaires, we derive trait-level scores reflecting the model’s implicit assumptions. Statistical and visualization analyses (Z-score normalization, ANOVA, PCA, and correlation mapping) are used to identify systematic differences across demographic groups.

The contributions of this paper are threefold:

1. It introduces a reproducible framework for psychometric elicitation from LLMs using established psychological instruments.
2. It performs a large-scale cross-demographic analysis, comparing Big Five and Dark Triad patterns across identities.
3. It offers interpretive insights into how implicit stereotype structures emerge in model-generated personas and discusses their ethical implications.

Through this approach, we aim to move beyond surface-level bias detection and reveal *how LLMs encode the psychology of stereotypes*—an essential step toward ensuring fairness, interpretability, and social responsibility in AI systems.

2. Related Work

The intersection of *bias analysis*, *psychometric evaluation*, and *Large Language Models (LLMs)* has become an emerging research domain, connecting machine learning with cognitive and social psychology. Existing literature largely focuses on linguistic, representational, or statistical bias — such as gendered associations in word embeddings, or disparities in model outputs across demographic identities. However, far fewer studies examine the psychological dimensions of these biases: how an LLM implicitly constructs the *personality* or *moral character* of different groups.

Recent advances in *persona-based prompting* have shown that LLMs can consistently simulate personality traits, preferences, and moral judgments when conditioned on contextual cues. This ability implies that underlying latent spaces in these models contain *consistent psychological mappings* learned from human discourse. Yet, those mappings may reflect — and potentially amplify — pre-existing cultural stereotypes present in the training data.

The present study builds upon this growing body of research by framing bias not merely as a statistical imbalance, but as a psychometric attribution phenomenon. In this view, an LLM’s response to personality-related prompts can be treated as a projection of internalized social constructs. This approach bridges three domains:

- LLM Bias Auditing,
- Computational Psychometrics, and
- Social Bias Theory in AI Ethics.

By situating our work within these areas, we extend previous studies that have analyzed bias at the textual and semantic level, moving toward a *cognitive-layer* interpretation of AI fairness.

2.1. Bias and Fairness in Large Language Models

The issue of bias in artificial intelligence has evolved from a technical concern into a central ethical challenge for AI research. In the context of Large Language Models (LLMs), *bias* refers to systematic and undesirable variations in model behavior that reflect or reinforce societal stereotypes, inequities, or cultural prejudices. Because LLMs are trained on massive text corpora collected from the internet, social media, and historical archives, they inevitably inherit the linguistic and cultural patterns present in those datasets. Studies have shown that this process leads to *encoded stereotypes* that manifest in model outputs — from gendered pronoun associations and occupational stereotypes to ideological bias in political or moral reasoning.

Fairness in LLMs is therefore a multifaceted concept. It encompasses:

- Representational fairness, i.e., ensuring that model embeddings do not encode discriminatory associations (e.g., “doctor” = male, “nurse” = female);
- Procedural fairness, ensuring equal performance across demographic subgroups.
- Outcome fairness, meaning that the model’s decisions or generated content do not disadvantage specific populations.

Research on bias mitigation in LLMs has included data filtering, controlled fine-tuning, reinforcement learning with human feedback (RLHF), and prompt-level interventions such as *debiasing templates* and *adversarial prompting*. However, most of these approaches treat bias as a *linguistic artifact*—an explicit surface-level phenomenon.

Recent work extends this perspective by examining latent bias: implicit patterns within the model’s internal representations that correspond to deeper social stereotypes. For example, certain demographic identifiers can shift the sentiment, tone, or emotional intensity of

responses, even when the semantic content remains neutral. Such findings suggest that LLMs encode *cognitive-like priors* about different demographic groups — a property that links bias to personality perception and social attribution mechanisms [1].

By situating fairness in a psychometric context, the current study explores a new question:

How does an LLM “imagine” the personality and moral traits of demographic identities?

This redefinition of fairness — from observable bias to *attributed bias* — enables a more granular understanding of how stereotype structures are generated within model cognition [2].

2.2. Psychometrics and Artificial Intelligence

Psychometrics — the quantitative study of psychological traits and personality — provides a rigorous framework for measuring latent dimensions of human cognition, emotion, and behavior. Over the past decades, personality models such as the Big Five and the Dark Triad have become standard instruments in both psychological research and computational modeling. Their structured, quantitative nature makes them ideal for integration with artificial intelligence systems seeking to emulate or analyze human-like behavior.

The Big Five Model, also known by the acronym EACNO (Extraversion, Agreeableness, Conscientiousness, Neuroticism, Openness), represents the most empirically validated taxonomy of personality.

- *Extraversion* captures sociability, assertiveness, and energetic engagement;
- *Agreeableness* reflects empathy, cooperation, and interpersonal warmth;
- *Conscientiousness* corresponds to organization, reliability, and self-discipline;
- *Neuroticism* denotes emotional instability and sensitivity to stress;
- *Openness to Experience* measures intellectual curiosity and creativity.

In contrast, the Dark Triad framework — consisting of *Machiavellianism* (M), *Narcissism* (NAR), and *Psychopathy* (PSY) — focuses on socially aversive traits that predict manipulative, exploitative, or self-serving tendencies. While these constructs often appear in psychological and criminological research, they have recently been adopted by computational social science to explore the moral and ethical dimensions of digital agents.

When applied to LLMs, these frameworks enable an unprecedented type of analysis: rather than evaluating model outputs purely for factual accuracy or bias, researchers can profile the model’s “personality” through its responses. Several studies have shown that GPT-type

models produce consistent Big Five profiles that can even vary with temperature settings or instruction style. This suggests that *latent personality structures* emerge from the statistical regularities of language learning itself.

Furthermore, mapping Dark Triad traits in LLM behavior reveals potential moral asymmetries — such as overconfidence, manipulateness, or emotional detachment — which mirror human dark-side cognition. Investigating these dimensions provides insight into the affective biases and moral priors encoded during model training.

By quantifying personality expression in LLM outputs, psychometric analysis serves as a diagnostic tool for *evaluating cognitive alignment* and *ethical safety*. It bridges the gap between surface-level text evaluation and deeper models of artificial “psychology.” In this study, psychometric scoring becomes the foundation for measuring how LLMs internalize demographic stereotypes — effectively translating social bias into measurable psychological variance [3],[4].

2.3. LLMs and Persona Conditioning

One of the most distinctive capabilities of modern Large Language Models (LLMs) lies in their contextual adaptability — the ability to modify style, tone, and reasoning according to the user’s prompt. This property, often referred to as persona conditioning, allows the model to adopt a specific identity, perspective, or emotional stance when instructed through natural language. For instance, prompting a model with “You are a compassionate therapist” or “You are a competitive entrepreneur” leads to consistent and thematically coherent response patterns.

This phenomenon has generated increasing academic interest, as it suggests that LLMs possess latent representation layers that encode human-like behavioral regularities. These representations can be activated or modulated through identity cues — including demographic descriptors such as gender, race, religion, or region. In other words, conditioning the model on an identity context effectively elicits the model’s internal stereotype of that persona.

Earlier works on persona simulation have shown that LLMs can maintain internal consistency across multiple responses, producing coherent personality profiles aligned with the given role. For example, when repeatedly asked Big Five or moral-dilemma questions, an LLM conditioned as a “female scientist” or a “religious leader” tends to generate reproducible psychometric signatures. Such consistency suggests that personas are not superficial textual masks, but stable attractors within the model’s conceptual space — emergent clusters of linguistic, emotional, and moral associations learned from training data.

From a psychological standpoint, persona conditioning parallels the process of stereotype activation in humans. When primed with demographic cues, individuals unconsciously draw on culturally learned scripts about how people from that group “think” or “behave.” Similarly, LLMs — having been trained on human-generated text — replicate these associative patterns in their outputs. The result is a computational form of implicit social cognition, in which the model reflects collective cultural expectations rather than neutral reasoning.

For researchers, this capability offers a double-edged tool. On one hand, it enables powerful simulations of social identities, useful for dialogue systems, storytelling, or empathy modeling. On the other, it exposes the internalized social biases of the model’s training distribution.

Therefore, analyzing LLM responses under controlled persona prompts provides an experimental gateway into understanding how language models reproduce demographic stereotypes — not through explicit prejudice, but through statistically learned personality and moral archetypes.

This study operates on persona conditioning as a systematic probing mechanism. By creating balanced combinations of gender, race, religion, and regional identity, and administering psychometric questionnaires to each synthetic persona, we can measure how the LLM’s attributed personality shifts across demographic dimensions. These controlled variations form the empirical backbone for identifying psychometric bias patterns in LLM-generated personas.

2.4. Research Gap

While the existing body of research on Large Language Model (LLM) bias has achieved significant progress in identifying linguistic disparities, it remains primarily constrained to surface-level phenomena—word associations, sentiment shifts, and topic preferences. These studies, although valuable, capture only the explicit layer of bias. They do not address how deeper cognitive-like structures within LLMs may encode *implicit psychological representations* of social groups.

Similarly, prior work on AI personality modeling has largely aimed at aligning machine behavior with human personality frameworks for interaction design or empathy generation. Few studies have examined personality attribution not as a *design feature*, but as a *diagnostic lens* for uncovering underlying biases.

While recent frameworks such as TRAIT [5] have successfully demonstrated that LLMs can maintain consistent personality profiles, they primarily focus on the

existence and consistency of these personas. Our work extends this methodology by repurposing psychometric instruments as a comparative fairness auditing tool. Rather than simply verifying that a model has a personality, we conduct a large-scale cross-persona and intersectional analysis to measure how that personality systematically degrades or shifts based on demographic attributes. This moves the utility of psychometrics from ‘persona design’ to ‘bias detection’. Most LLM personality studies assume a single, “universal” model personality rather than exploring how that personality fluctuates when the model is prompted with diverse demographic identities.

Furthermore, the Dark Triad dimension — representing Machiavellianism, Narcissism, and Psychopathy — has been almost entirely absent from fairness and bias research in artificial intelligence. These traits, although negatively connoted, provide crucial insight into *moral asymmetries* and *affective biases*. Understanding how LLMs distribute these traits across demographics can reveal implicit associations between identity and morality encoded in training data.

Another methodological gap concerns cross-dimensional bias interaction. Most evaluations focus on single-axis demographics (e.g., only gender or only race). In contrast, real-world stereotypes are *intersectional*, emerging from combinations such as “female–religious–Asian” or “male–atheist–Western European.” This study addresses that limitation by systematically varying four demographic factors — gender, race, religion, and region — across a large, balanced persona set.

Finally, while recent bias audits use quantitative fairness metrics, they often lack interpretability. Traditional bias measures (e.g., KL divergence or accuracy gaps) reveal *that* differences exist but not *how* they manifest semantically or psychologically. By applying psychometric frameworks (Big Five and Dark Triad) to LLM outputs, this study introduces a human-interpretable metric of bias, translating abstract probability shifts into personality trait differences.

In summary, the key research gaps this work addresses are:

1. From surface bias to latent bias: Moving beyond textual stereotypes to cognitive-level psychometric associations.
2. From general personality to differential attribution: Measuring how LLMs alter personality traits across demographic identities.
3. From fairness metrics to interpretability: Using established psychological taxonomies to explain *how* and *why* demographic stereotypes emerge.
4. From single axis to intersectional analysis: Exploring multi-factor demographic bias patterns.

By filling these gaps, this research contributes a novel interdisciplinary framework that merges computational linguistics, psychometrics, and AI ethics — advancing the discussion of fairness in LLMs toward the domain of *machine social cognition* [6].

3. Methodology

3.1. Persona Generation Framework

To investigate how Large Language Models (LLMs) implicitly encode demographic stereotypes through psychometric attributions, we developed a structured persona generation framework. This framework systematically combines demographic categories to create balanced and reproducible *synthetic identities* that can be used to probe model behavior.

Each persona is defined across four demographic dimensions — *region*, *gender*, *race*, and *religion* — producing a diverse set of cultural and social contexts. The following categories were used:

- Geopolitical Regions (11 total):
Western Europe, Eastern Europe, North America, Latin America, Middle East, Sub-Saharan Africa, South Asia, East Asia, Southeast Asia, Central Asia, and Oceania.
- Races (5 total):
White, Black, Asian, Latino, and Mixed.
- Religions (6 total):
Orthodox Christian, Catholic, Muslim, Buddhist, Hindu, and Atheist.
- Genders (2 total):
male and female.

The full factorial combination of these categories' yields:

$$11 \text{ regions} \times 5 \text{ races} \times 6 \text{ religions} \times 2 \text{ genders} = 660 \text{ unique personas.}$$

Each persona represents a unique demographic identity prompt. To generate responses, every persona was presented to the model using a standardized prompt template:

"You are a {gender}, {race}, {religion} average person from {region}.

Answer the following question as such a person would respond on a scale from 1 to 5

(1 = Strongly Disagree, 5 = Strongly Agree):"

This template was selected for its clarity, neutrality, and balanced linguistic framing. By introducing demographic identity markers without evaluative or emotional language, it encourages the LLM to generate responses based on *implicit cultural priors* rather than explicit instructions. Each persona was queried sequentially across a full battery of psychometric items (50 for the Big Five and 12 for the Dark Triad). For every (persona, question) pair, the model produced a

numerical Likert response (1–5), which was stored in structured form along with question metadata. The resulting dataset was composed of:

- 660 personas,
- 62 questions per persona,
- yielding a total of 40,920 recorded responses.

Figure 1 below summarizes and corroborates the experimental design detailed above, visualizing the workflow from the full factorial combination of demographic attributes to the generation of 660 unique personas and the subsequent collection of 40,920 quantitative responses.

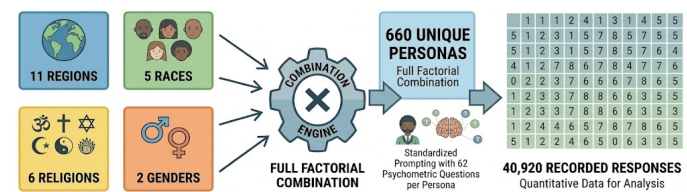


Figure 1: Descriptive Overview of the Psychometric AI Persona Study Data Generation Pipeline.

Data collection was performed automatically using Python, with deterministic decoding to ensure reproducibility. The persona generation loop iterated through all category combinations, formatted the prompts, queried the model, and stored responses in a unified dataframe (persona_results). A simplified version of the procedure is shown below:

This process effectively transforms the LLM into both a *subject* (producing the responses) and an *object of study* (whose internal biases are measured). Each persona acts as a controlled probe, enabling cross-demographic comparison of the model's psychometric attributions.

The output of this framework is a structured dataset — df_full — containing all persona identities, questions, and Likert-scale answers. This dataset constitutes the empirical foundation for all subsequent analyses described in Sections 3.2–3.6 [7],[8].

3.2. Questionnaire Design

The psychometric questionnaire used in this study was designed to elicit *structured personality responses* from the LLM across two major theoretical frameworks: (1) the Big Five Personality Model (EACNO), and (2) the Dark Triad Model (SD3). Together, these frameworks capture both prosocial and antisocial personality dimensions, providing a comprehensive basis for evaluating how the model attributes character traits to different demographic personas [9].

We adopted a standardized questionnaire approach similar to established datasets like TRAIT [5]; however, we

significantly expanded the scope of evaluation. Instead of testing for internal consistency within a single persona, our framework applies these instruments across a full factorial combination of 660 demographic identities. This allows us to isolate specific attribute-based distortions (e.g., how changing only 'religion' alters perceived 'conscientiousness'), effectively turning the questionnaire into a differential diagnostic for latent stereotypes.

3.2.1. Big Five Personality Items

The Big Five Model represents the gold standard of personality psychology, quantifying personality along with five independent factors: Extraversion (E), Agreeableness (A), Conscientiousness (C), Neuroticism (N), and Openness to Experience (O).

A set of 50 Likert-scale statements was employed to evaluate these five traits (10 items per trait). The items were adapted from validated short-form Big Five inventories (e.g., the International Personality Item Pool – IPIP) and rephrased for clarity and simplicity to suit LLM prompting. Each item expresses a self-assessment statement such as:

"I see myself as someone who is talkative."
"I get chores done right away."
"I worry a lot."
"I am original and come up with new ideas."

To maintain psychometric integrity, reverse-coded items were preserved where applicable. For example, low Extraversion items such as *"I am reserved"* were included and scored inversely during post-processing. This balance prevents the model from simply pattern-matching affirmative phrasing and ensures that the variance of responses reflects underlying psychological consistency. Each of the 50 items was presented as a separate prompt within the persona context. The model's numeric response (1–5) to each item was stored as `best_answer`, corresponding to the following [9].

Likert Structure:

1. Strongly Disagree
2. Disagree
3. Neutral
4. Agree
5. Strongly Agree

3.2.2. Dark Triad (SD3) Items

To complement the Big Five, we incorporated 12 items derived from the *Short Dark Triad (SD3)* instrument (Jones & Paulhus, 2014), covering three subscales:

- **Machiavellianism (M)** — manipulativeness, strategic deception, and pragmatic morality

- **Narcissism (NAR)** — grandiosity, self-focus, and need for admiration
- **Psychopathy (PSY)** — impulsivity, callousness, and emotional detachment

Each subscale was assessed through four statements. Example prompts included:

"I manipulate others to get my way."
"I insist on getting the respect I deserve."
"I lack remorse after hurting someone."

As with the Big Five, the same 1–5 Likert scale was used, ensuring consistency across the psychometric space.

The inclusion of Dark Triad traits extends the analysis beyond classical personality constructs, enabling the study of moral asymmetry in model behavior — i.e., whether the LLM assigns morally "darker" traits more frequently to certain demographics [9].

3.2.3. Adaptation for LLM Context

Unlike human participants, LLMs do not possess self-awareness or emotions. Therefore, the questionnaire was restructured to simulate *third-person perspective attribution*: the prompts instructed the model to respond as if it were the average person from a given demographic group, rather than as itself. This reframing allowed the model to project *collective cultural knowledge* rather than introspection [9].

Each prompt explicitly stated:

"Answer the following question as such a person would respond..."

This phrasing reduces the likelihood of meta-cognitive replies (e.g., *"As an AI language model, I cannot feel emotions"*) and constrains the model within a behavioral simulation space. Pilot tests confirmed that this phrasing yielded stable numeric outputs across multiple runs, indicating consistent interpretation.

To verify psychometric coherence, inter-item correlations were examined post hoc, and the response patterns exhibited meaningful variance across traits and demographics — validating the use of the adapted questionnaire as a diagnostic probe for LLM stereotypes.

3.3. Trait Computation and Scoring

Following data collection, each persona's responses were aggregated into numerical trait scores according to standardized psychometric scoring procedures. The scoring framework combined established Big Five (EACNO) and Dark Triad (SD3) computation schemes, adapted for automated calculation within the experimental pipeline [5].

3.3.1. Big Five (EACNO) Scoring

The Big Five personality traits were computed based on the scoring scheme of the International Personality Item Pool (IPIP) short-form inventory, using 10 items per trait.

For each trait, positive and reverse-coded items were weighted accordingly to preserve scale directionality. The raw scores were calculated as follows:

$$\begin{aligned} E &= 20 + Q_1 - Q_6 + Q_{11} - Q_{16} + Q_{21} - Q_{26} + Q_{31} - Q_{36} + Q_{41} - Q_{46} \\ A &= 14 - Q_2 + Q_7 - Q_{12} + Q_{17} - Q_{22} + Q_{27} - Q_{32} + Q_{37} + Q_{42} + Q_{47} \\ C &= 14 + Q_3 - Q_8 + Q_{13} - Q_{18} + Q_{23} - Q_{28} + Q_{33} - Q_{38} + Q_{43} + Q_{48} \\ N &= 38 - Q_4 + Q_9 - Q_{14} + Q_{19} - Q_{24} - Q_{29} - Q_{34} - Q_{39} - Q_{44} - Q_{49} \\ O &= 8 + Q_5 - Q_{10} + Q_{15} - Q_{20} + Q_{25} - Q_{30} + Q_{35} + Q_{40} + Q_{45} + Q_{50} \end{aligned}$$

where Q_i denotes the Likert score (1–5) for question i . Positive and negative signs represent normal or reverse-coded items respectively. The additive constants (e.g., 20, 14, 38, 8) ensure that the resulting values fall within interpretable personality scale ranges consistent with the IPIP framework.

Each computed value corresponds to a **trait magnitude** per persona, expressing the LLM's inferred intensity of that characteristic when role-playing as a member of the corresponding demographic group.

To verify internal consistency, the resulting distributions were examined for:

- variance across personas (ensuring diversity of LLM attributions),
- and inter-trait correlation patterns (confirming expected psychological relationships, e.g., E positively correlated with O and negatively with N) [5].

3.3.2. Dark Triad (SD3) Scoring

The Short Dark Triad (SD3) instrument was used to quantify the model's attribution of socially aversive or morally self-centered traits. Each of the three Dark Triad dimensions — *Machiavellianism* (M), *Narcissism* (NAR), and *Psychopathy* (PSY) — was computed as the sum of four corresponding items:

$$\begin{aligned} M &= Q_{51} + Q_{52} + Q_{53} + Q_{54} \\ NAR &= Q_{55} + Q_{56} + Q_{57} + Q_{58} \\ PSY &= Q_{59} + Q_{60} + Q_{61} + Q_{62} \end{aligned}$$

The resulting values represent each persona's estimated "dark trait intensity", derived from the model's Likert-scale responses. Because the range of each item is 1–5, each Dark Triad subscore spans 4–20. Larger scores indicate stronger endorsement of manipulative, egocentric, or emotionally detached tendencies [5].

3.3.3. Automation and Validation

All computations were executed programmatically in Python to ensure repeatability and minimize human bias. Each persona's response vector (62 items) was indexed by

question_id and processed through automated formulas that replicated the IPIP and SD3 scoring structure.

Each persona's results were stored in a consolidated dataframe (`df_scores`) with eight columns: ' $E, A, C, N, O, M, NAR, PSY$ '.

Descriptive analysis confirmed logical consistency:

- E (Extraversion) and NAR (Narcissism) showed moderate positive correlation,
- A (Agreeableness) negatively correlated with M (Machiavellianism) and PSY (Psychopathy), reflecting realistic psychological interdependencies — a strong indicator that the LLM internalized culturally plausible personality structures [5].

3.4. Data Normalization and Z-Scoring

Before performing any comparative or inferential analysis, it was essential to normalize the computed personality and Dark Triad scores to a common scale. Raw scores derived from the Big Five and SD3 inventories differ in their numerical range and variance: for example, *Extraversion* values typically span 10 – 50, whereas *Machiavellianism* ranges only 4 – 20. Directly comparing such values could therefore exaggerate or obscure cross-trait differences. To address this issue, all scores were standardized using Z-score normalization.

3.4.1. Z-Score Formula

For each trait $t \in \{E, A, C, N, O, M, NAR, PSY\}$, the Z-score for persona i was computed as:

$$Z_{i,t} = \frac{X_{i,t} - \mu_t}{\sigma_t}$$

where

- $X_{i,t}$ is the raw trait score for persona i ,
- μ_t is the mean score of trait t across all personas, and
- σ_t is the standard deviation of trait t across all personas.

This transformation centers each trait around zero mean and unit variance, producing dimensionless values that are directly comparable across both traits and demographic groups.

In practice, positive Z-values indicate that a persona scores above the global average for a given trait, whereas negative values indicate below-average representation. This allows for an intuitive interpretation of bias: a consistent positive deviation for a demographic group suggests a systematic over-attribution of that trait by the model.

3.4.2. Implementation

The resulting standardized dataset (`df_scores_z`) preserved the original persona identifiers while replacing raw trait values with Z-scores.

Each persona thus corresponds to an eight-dimensional normalized feature vector, enabling cross-group statistical comparison.

3.4.3. Analytical Use

The normalized dataset served as the foundation for all subsequent statistical and visualization analyses, including:

- Heatmaps of mean Z-scores per demographic group (Figures 1–2) to visualize bias direction and magnitude.
- Bar and radar plots, highlighting which personas or groups were most atypical relative to the overall population mean.
- ANOVA and t-tests, applied to standardized scores to detect significant group-level differences without scale distortion.
- Principal Component Analysis (PCA), leveraging the zero-mean normalization to identify latent clusters in trait space.

Z-score normalization not only ensured mathematical comparability but also enabled psychological interpretability: each deviation of one standard deviation represents a meaningful difference in trait attribution strength, facilitating a consistent interpretation of bias magnitude across all dimensions.

3.5. Statistical Analysis and Visualization

Once the psychometric and Dark Triad scores were computed and normalized, a series of statistical and visualization techniques were applied to quantify demographic bias and reveal latent personality structures within the LLM's responses. The analysis was designed to examine both *group-level differences* and *underlying correlations* between traits, providing complementary perspectives on model behavior.

3.5.1. Group-Level Analysis (ANOVA and t-tests)

To determine whether the LLM assigned significantly different personality or moral traits to different demographic categories, we performed Analysis of Variance (ANOVA) tests for each trait across the four main demographic factors: *gender*, *race*, *religion*, and *region*.

For each trait t , the one-way ANOVA model was defined as:

$$H_0: \mu_{1t} = \mu_{2t} = \dots = \mu_{kt} \text{ vs. } H_a: \text{at least one group mean differs.}$$

Here, μ_{jt} represents the mean Z-score of trait t within group j (e.g., male vs. female). A statistically significant p -value ($p < 0.05$) indicates that the model exhibits systematic differentiation in how it assigns that trait across demographic groups.

Following ANOVA, pairwise Welch t-tests were conducted to identify which specific groups differed. These pairwise comparisons yielded two key outputs:

- Mean difference (Δ), representing the direction and magnitude of bias; and
- p -value, quantifying statistical significance.

For example, if *Agreeableness* (A) showed $\Delta = -0.45$ (female–male) and $p = 0.02$, this was interpreted as the model attributing higher *Agreeableness* to female personas.

This analysis produced a structured bias matrix per factor, later visualized as heatmaps and bar charts (Figure 1C, Tables 1–2).

3.5.2. Correlation Analysis

To explore inter-trait dependencies and psychometric coherence, a correlation matrix was computed across all eight dimensions (E, A, C, N, O, M, NAR, PSY). The Pearson correlation coefficient r was used to quantify the linear relationships between traits:

$$r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

The resulting correlation heatmap (Figure 4) revealed patterns consistent with psychological theory — for instance, strong negative correlation between *Agreeableness* and *Psychopathy* ($r \approx -0.6$), and positive correlation between *Extraversion* and *Narcissism* ($r \approx +0.4$). Such patterns support the interpretive validity of the LLM's simulated personalities and confirm that the model expresses *internally consistent personality structures*, not random noise.

3.5.3. Principal Component Analysis (PCA)

To visualize the overall structure of LLM-generated personas, Principal Component Analysis (PCA) was applied to the Z-score matrix. This unsupervised dimensionality reduction technique identifies orthogonal components that capture the greatest variance in the dataset:

$$Z = W \cdot P$$

where W represents the component weights and P the principal component loadings.

The first two principal components (PC1, PC2) explained approximately 60–70% of the total variance, forming a two-dimensional *trait map*. Personas were then plotted in this reduced space, colored by demographic attributes (e.g., race, region, gender). Distinct clustering patterns (Figure 3) indicated that certain groups shared similar psychometric profiles — evidence of consistent stereotype formation within the model's latent space.

Outliers identified in the PCA corresponded to demographic combinations that the model associated with

particularly extreme trait attributions (e.g., high Narcissism or low Agreeableness). These clusters were interpreted as *bias attractors*, representing the LLM’s internalized archetypes.

3.5.4. Visualization Framework

To communicate effectively, several complementary visual representations were generated using Python libraries such as matplotlib and seaborn:

- Heatmaps: visualized group-level Z-score averages, highlighting direction and magnitude of demographic bias.
- Boxplots: displayed raw score distributions per demographic category to show score dispersion and overlap.
- Bar charts: ranked differences (Δ) in trait attribution (e.g., male vs. female).
- Radar charts: compared normalized profiles across top 3 most divergent groups (e.g., races or regions).
- PCA scatter plots: visualized latent psychometric clusters.
- Correlation maps: revealed structural relationships between traits.

Each visualization was exported in high-resolution PNG format and labeled according to the JENRS figure standard (Figures 1–4). Together, these figures constitute an interpretable visual narrative of how the model’s internal representation space mirrors human social cognition and bias.

3.5.5. Summary of Statistical Pipeline

The complete analytical workflow is summarized as follows in Table 1.

Table 1: Summary of Statistical Pipeline

Step	Method	Purpose
1	One-way ANOVA	Test group-level differences per trait
2	Pairwise t-tests	Identify directionality and strength of bias
3	Z-score normalization	Standardize scale across traits
4	PCA	Visualize latent personality clusters
5	Correlation matrix	Verify psychometric coherence
6	Visualization	Present interpretable findings

This integrated approach allows both quantitative rigor and qualitative interpretability, bridging computational bias detection with psychological insight.

3.6. Technical Implementation Environment

All data collection, trait computation, and statistical analyses were implemented in Python, using a fully reproducible software environment. The computational pipeline was designed to ensure transparency, replicability, and scalability across different LLM configurations.

3.6.1. Software Framework

The entire workflow — from persona generation to statistical visualization — was implemented as a modular Python project. The following libraries were employed as shown in Table 2:

Table 2: Libraries Table

Library	Purpose
Pandas	Data manipulation, tabular storage of responses (df_full, df_scores, df_scores_z)
Numpy	Numerical computation and array operations
scipy.stats	Statistical analysis, Z-score normalization, t-tests, and ANOVA
matplotlib / seaborn	Visualization (heatmaps, barplots, radar charts, PCA scatterplots)
scikit-learn	Dimensionality reduction via PCA
Openpyxl	Exporting structured results to Excel format
Tqdm	Progress tracking during persona generation
transformers / huggingface_hub	Interfacing with the selected LLM model
random / itertools	Deterministic iteration through demographic combinations

The modularity of the framework allows each component — prompt generation, response collection, scoring, and visualization — to operate independently while sharing a common data schema.

3.6.2. Model and Prompt Execution

All responses were obtained from a Large Language Model (LLM) using deterministic inference parameters to ensure experimental consistency.

The model used in this study was Meta LLaMA-3.1-8B-Instruct, deployed via the Hugging Face Transformers API.

Inference parameters:

- Temperature: 0.0 (deterministic sampling)
- Top-p (nucleus sampling): 1.0
- Max tokens: 256
- Repetition penalty: 1.0
- Stop sequences: newline and “Answer:” markers

Each prompt followed the structured format described in Section 3.1. The use of deterministic decoding (temperature = 0) ensured that identical personas and questions always yielded identical responses, enabling one-to-one comparison across demographic groups.

Response parsing and token probability extraction were automated using a custom wrapper function `get_token_probs()`, which computed the likelihood of each Likert-scale response (1–5) and selected the one with the highest probability as the model’s “answer.”

The model used in this study was Meta LLaMA-3.1-8B-Instruct, deployed via the Hugging Face Transformers API. The primary experiments were conducted using LLaMA-3.1-8B-Instruct due to its open-weight availability, strong instruction-following performance, and widespread adoption in recent LLM research. This model provides an appropriate balance between representational capacity and experimental reproducibility, making it suitable for systematic bias analysis.

3.6.3. Computational Environment

All experiments were conducted on a high-performance local mobile workstation with the following specifications as shown in Table 3:

Table 3: Local mobile workstation specifications

Component	Specification
CPU	AMD Ryzen 7 8845HS (8 cores / 16 threads)
RAM	48 GB DDR5
GPU	NVIDIA RTX 4060 (8 GB VRAM)
Storage	2 TB NVMe SSD
Operating System	Windows 11 Pro (64-bit)
Python Version	3.11
CUDA-Support	Enabled via Transformers

The model weights and tokenizer were loaded locally to minimize latency and ensure complete control over inference settings. All intermediate results, figures, and tables were saved under versioned directories (e.g., `/report_export/`, `/final_figures/`) for reproducibility.

3.6.4. Reproducibility and Version Control

To guarantee reproducibility, random seeds were fixed across all scripts, and the same persona order was maintained during every experimental run. Version control was managed through **Git**, ensuring that code, data, and results could be tracked and replicated. Additionally, all generated Excel outputs (e.g., `persona_answers_scores_with_zscores.xlsx`) were timestamped and stored with metadata (model version, date, system hash).

This technical architecture ensures that any researcher can replicate the study by:

1. Running the provided Python scripts,
2. Supplying the same demographic combinations and questionnaire items, and
3. Using an equivalent LLM configuration.

3.6.5. Workflow Summary

The full experimental workflow can be summarized as:

1. Persona Definition → generation of demographic combinations
2. Prompt Execution → querying the LLM with psychometric items
3. Response Parsing → extracting Likert-scale outputs
4. Trait Scoring → computing EACNO and SD3 dimensions
5. Normalization → applying Z-score transformation
6. Statistical Testing → ANOVA, t-tests, correlation, PCA
7. Visualization → generating figures and summary heatmaps
8. Reporting → exporting Excel sheets and publication-ready figures

This pipeline integrates both *psychological modeling* and *computational reproducibility*, forming a robust foundation for demographic stereotype elicitation in LLMs.

Figure 2 below illustrates the end-to-end experimental workflow, integrating the entire pipeline into five distinct stages. The process advances from Persona Construction and Prompting to the generation of LLM Responses, which are subsequently quantified during Scoring and evaluated in the final Analysis phase.

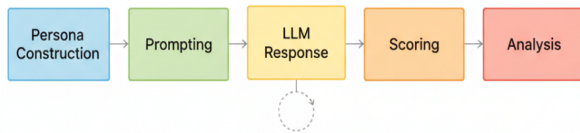


Figure 2: LLM Experimentation workflow.

4. Results

The LLM-generated personas exhibit distinct trait patterns across different demographic categories. As an initial overview, as we can see in Figure 3 (panels A–C) summarizes the mean standardized trait scores (Z-scores) for each demographic group in race, religion, and region, while panel D provides a radar chart comparing the multi-trait profiles of three illustrative racial groups. In these heatmaps, pronounced color differences immediately suggest stereotype-consistent biases. For example, panel A highlights that personas with Mixed race have starkly higher scores on dark traits (deep red in columns M, NAR, PSY) coupled with much lower Big Five scores (deep blue in E, A, C), whereas other races show more moderate hues. Panel B suggests that Atheist personas (top row) diverge strongly on certain traits (notably dark blue for A and C indicating very low Agreeableness and Conscientiousness). Panel C focuses on a subset of regions with the largest deviations, revealing, for instance, North America's lower Machiavellianism (blue in column M) and Oceania's higher Neuroticism (red in N). The radar chart in panel D further illustrates how an entire trait profile can differ by race: the Mixed profile (blue shaded area) bulges out dramatically along the dark triad axes compared to the Latino (orange) and Black (green) profiles, which extend more on positive personality trait axes. The following subsections provide a detailed breakdown of these patterns for each demographic dimension.

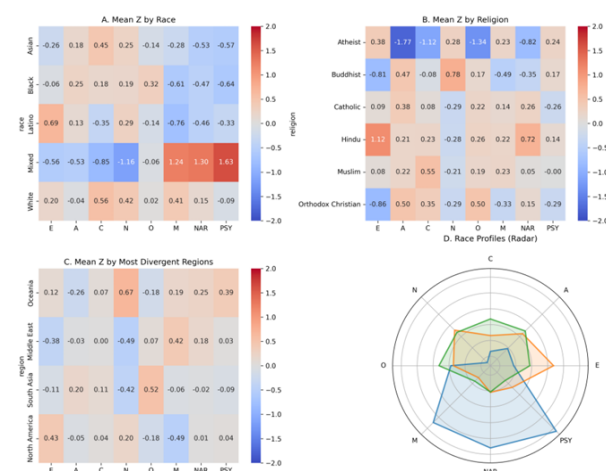


Figure 3: Overview of demographic biases in trait scores. Panel A – Mean Z-scores by Race; Panel B Mean Z-scores by Religion; Panel C – Mean Z-scores by Region Panel D – Radar chart of trait profiles for select races (Mixed, Latino, Black).

4.1. Regional Trait Differences

Regional origin is associated with systematic variations in persona trait profiles as shown in Figure 4. Clear patterns emerge in the Big Five dimensions across regions. Extraversion (E) tends to be highest for Western, English-speaking regions (e.g., Western Europe and North America) and lowest for regions like Central Asia and the Middle East, indicating a stereotype of Western personas as more outgoing and certain Asian/Middle Eastern personas as more introverted. Agreeableness (A) varies less extremely, but Central Asia stands out with a notably low A (a stereotype of lower cooperativeness) while regions such as South Asia and Latin America are slightly higher than average. Conscientiousness (C) is depicted as relatively high in parts of Asia (e.g., Southeast Asia) and lower in some Western or African regions (e.g., Western Europe and Sub-Saharan Africa). Neuroticism (N) shows one of the widest gaps: Oceania has a very high average N (suggesting personas from Oceania are portrayed as especially prone to anxiety), whereas the Middle East and Eastern Europe have very low N (stereotyping those personas as emotionally stable or stoic). Openness (O) also differs by region: South Asia is highest (implying very open-minded personas), whereas East Asia is lowest, with Central Asia and Oceania also somewhat lower (indicating more traditional or less open portrayals for those regions).

Turning to the Dark Triad traits, we see distinctive regional stereotypes as well. Machiavellianism (M) is notably high for Middle Eastern personas (the only region markedly above average) and lowest for North American personas, suggesting that the model tends to cast Middle Eastern characters as more manipulative and North American characters as more straightforward. Most other regions hover near the average on M (lighter colors), with slight positive bias in some (e.g. Southeast Asia) and slight negative in others (e.g. Western Europe). Narcissism (NAR) varies only slightly by region; no group deviates far from the mean (all around ± 0.2 Z). The Middle East and Latin America show mildly elevated NAR, whereas Western Europe is a bit below average, indicating only minor shifts in self-centeredness across locales. Psychopathy (PSY) has moderate regional differences: Oceania shows a higher PSY than most regions, and Latin America also has a modest elevation, meaning personas from these regions are depicted as somewhat more impulsive or low empathy. In contrast, Eastern and Western Europe have the lowest PSY (personas portrayed as more empathetic and rule-abiding). In summary, regional stereotypes in the model's outputs manifest as distinct personality profiles: for example, Western Europe and North America come across as more extraverted and conscientious but less Machiavellian; Central Asia and the Middle East as more introverted (and, in the Middle East's

case, more manipulative but less neurotic); and Oceania as notably more neurotic (and slightly more psychopathic) relative to others [10].

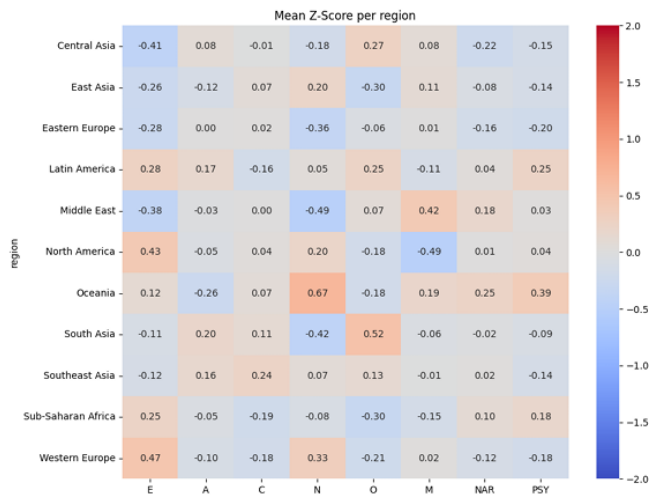


Figure 4: Mean Z-score per region. Heatmap of average standardized trait scores for personas from 11 global regions.

4.2. Religious Bias Patterns

Religious affiliation of the persona corresponds to strong divergences in the attributed traits as shown in Figure 5. Perhaps the most striking pattern is seen with Atheist personas, which deviate dramatically from all religious groups on multiple traits. Atheist profiles are characterized by very low Agreeableness ($A \approx -1.77$) and Conscientiousness ($C \approx -1.12$) — shown as dark blue cells — indicating that non-religious personas were overwhelmingly portrayed as less warm/compassionate and less dutiful/organized. They also show a notably low Openness ($O \approx -1.34$), suggesting a stereotype of close-mindedness or conventionality in atheist personas. These values are far below those of any religious group; for comparison, the next lowest Openness among religious categories is Orthodox Christian at -0.50 , and no religious group comes close to the extreme negative Agreeableness of the atheist group. Atheist personas further have moderately elevated dark traits: Machiavellianism ($M = +0.23$) and Psychopathy ($PSY = +0.24$) are slightly above average for atheists, whereas most religious groups hover around zero or below on these traits. Their Narcissism ($NAR = -0.82$) is lower than average, implying that despite being depicted as disagreeable, atheist personas are not shown as particularly narcissistic (if anything, somewhat humble or self-effacing, given the negative z-score).

In contrast, personas with religious identities generally cluster closer to the population's mean on most traits, with a few notable biases for each religion. Hindu personas stand out for exceptionally high Extraversion ($E \approx +1.12$, the reddest cell in column E) — depicting Hindu individuals as especially sociable or outgoing. Hindu profiles also show a pronounced spike in Narcissism ($NAR \approx +0.72$, bright red), making them the most

narcissistic on average among the groups. Other traits for Hindus are moderately above average ($A \approx +0.21$, $C \approx +0.23$, $O \approx +0.26$) with no strong negatives, meaning the LLM tended to imbue Hindu personas with generally positive Big-Five traits alongside the high extraversion and narcissism. Muslim personas, meanwhile, are characterized by the highest Conscientiousness ($C \approx +0.55$) among the religions — a substantial positive deviation (shown in red) suggesting a stereotype of Muslims as especially disciplined or responsible. Muslims also have slightly above-average Agreeableness and Openness ($A \approx +0.22$, $O \approx +0.19$) and near-average Extraversion ($E \approx +0.08$). Their dark trait scores are unremarkable: Machiavellianism is mild ($+0.23$, similar to Atheists), Narcissism about average ($+0.05$), and Psychopathy essentially zero, indicating no strong dark trait bias for Muslim personas aside from a minor Machiavellian lean.

Two groups, Buddhist and Orthodox Christian personas, both exhibit high Agreeableness ($A \approx +0.47$ and $+0.50$, respectively), marking them as the most agreeable (warm and cooperative) profiles among the set. They differ, however, in other traits. Orthodox Christian personas have very low Extraversion ($E \approx -0.86$, deep blue), meaning they are depicted as far more introverted or reserved. They also have moderately high Conscientiousness ($C \approx +0.35$) and markedly low Machiavellianism ($M \approx -0.33$) and Psychopathy ($PSY \approx -0.29$). This paints a stereotype of Orthodox Christian individuals as kind, dutiful, and non-manipulative — a generally prosocial profile. Buddhist personas, on the other hand, also show low Extraversion ($E \approx -0.81$) but combine it with one of the highest Neuroticism scores ($N \approx +0.78$) among the groups, suggesting a portrayal of Buddhists as relatively anxious or emotionally reactive despite being agreeable. Interestingly, Buddhists have the lowest Machiavellianism of all ($M \approx -0.49$, a dark blue cell in column M), aligning with a stereotype of high altruism or straightforwardness. Their Narcissism is slightly below average ($NAR \approx 0.35$) and Psychopathy slightly above average ($PSY \approx +0.17$). The combination for Buddhists is thus: modest, kind, somewhat anxious, and non-manipulative, with a hint of impulsivity (higher psychopathy) — a nuanced mix likely reflecting specific narrative tropes.

Catholic personas do not display extreme outliers on most traits; they remain closer to the population mean (mostly neutral-colored cells). They show a mildly higher Agreeableness ($A \approx +0.38$) comparable to the other religious groups and a slightly elevated Narcissism ($NAR \approx +0.26$). Notably, Catholics share a trend with Orthodox Christians of lower Psychopathy ($PSY \approx -0.26$ for Catholics, similar to Orthodox's 0.29), indicating that Christian-affiliated personas (both Catholic and Orthodox) were depicted as less psychopathic (more empathetic or rule-abiding). Catholics' Extraversion,

Conscientiousness, and Machiavellianism are all near zero ($E \approx +0.09$, $C \approx +0.08$, $M \approx +0.14$), suggesting no strong stereotype on those dimensions beyond general sociability and decency.

In summary, the LLM's personas reflect distinct religious stereotypes in trait attributes. Non-religious (Atheist) characters are cast in a particularly negative light on key prosocial traits (agreeableness, conscientiousness, openness) and somewhat higher in callousness-related traits, whereas each religious group carries its own subtle bias: Hindus as outgoing and narcissistic, Muslims as dutiful and reasonably well-rounded, Buddhists as kind yet anxious and least manipulative, Orthodox Christians as introverted, kind, and law-abiding, and Catholics as generally average with slight leanings toward kindness and low psychopathy. These findings suggest that rather than functioning as neutral arbiters, LLMs may inadvertently reinforce deep-seated societal prejudices. Consequently, the deployment of such models risks perpetuating historical tropes, potentially marginalizing specific groups through automated, biased characterizations [9].

4.3. Racial Trait Attribution

Significant trait biases are evident across different racial categories as shown in Figure 6. The most pronounced pattern is observed for the Mixed-race personas, who emerge as extreme outliers in the dataset. Mixed-race personas are portrayed with dramatically negative Big Five traits alongside highly elevated Dark Triad traits. In fact, they exhibit the lowest Extraversion, Agreeableness, and Conscientiousness of all races (far below the mean in those traits), suggesting a stereotype of Mixed individuals as especially unsociable, uncooperative, and undisciplined. At the same time, the Mixed group has by far the highest Machiavellianism, Narcissism, and Psychopathy scores, implying that when the persona's race is "Mixed," the model often imbues the character with an antagonistic, anti-social personality profile (manipulative, self-centered, and callous). This extreme combination – low Big Five coupled with high Dark Triad – is unique to the Mixed group in the model's output.

Other racial groups have more moderate, often favorable profiles. Latino personas, for example, are characterized by relatively positive social traits. They have the highest Extraversion of any race (indicating Latino characters are frequently depicted as very outgoing and energetic), and their Dark Triad scores are notably low. Machiavellianism for Latinos is extremely low (suggesting a stereotype of Latinos as very non-manipulative or straightforward), and both Narcissism and Psychopathy are below average as well. Latinos' Agreeableness and

Openness are roughly average (no strong bias), and Conscientiousness is slightly below average. Overall, the LLM portrays Latino personas as sociable and generally friendly, with a clear absence of "dark" characteristics – a stark contrast to the Mixed-race profile. Black personas similarly skew toward favorable Big Five attributes and low dark traits. They have the highest Agreeableness and Openness among the races, implying Black individuals are often depicted as particularly friendly, cooperative, and open-minded. Their Conscientiousness is also modestly above average. Importantly, Black personas have uniformly low Dark Triad scores: Machiavellianism, Narcissism, and Psychopathy are all significantly below zero, indicating a consistent tendency for the model to depict Black characters as less manipulative, less self-absorbed, and less psychopathic relative to the norm. Their Extraversion is about neutral. This trait pattern – high A and O coupled with low M/NAR/PSY – suggests an overall stereotype of Black personas as affable, well-adjusted, and trustworthy.

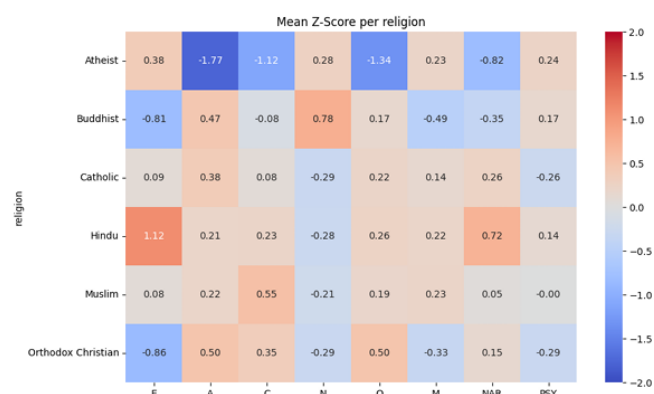


Figure 5: Mean Z-score per religion.

Asian personas have a distinct but comparatively balanced profile. They are depicted as more conscientious than others (C is relatively high, second only to White) and somewhat more agreeable than average. However, Asian characters tend to be shown as more introverted (low E) and a bit less open (slightly low O) in the model's outputs. In terms of dark traits, Asian personas are assigned uniformly low values: low Narcissism and Psychopathy, along with moderately low Machiavellianism. These indicate that Asian characters are stereotyped as polite, diligent, and non-antisocial – essentially a reserved but well-intentioned profile. They lack the strong sociability of the Latino group or the high openness of the Black group but also avoid any hint of the antagonistic Dark Triad elevation seen in Mixed personas. White personas tend to be portrayed near the average on most traits, with a couple of mild leanings. They have the highest Conscientiousness of all races, suggesting a stereotype of White individuals as especially organized or responsible. Their Extraversion is slightly above the mean as well (though not as high as Latinos), and Neuroticism is somewhat elevated (indicating White personas might be depicted as a bit more

prone to stress or negative emotions compared to others). White personas' Machiavellianism is mildly above average (the highest after Mixed-race, though far below the extreme Mixed value), implying a small bias toward portraying White characters as somewhat more strategic or manipulative than most other groups. Their Narcissism is also slightly positive and Psychopathy slightly negative (effectively near neutral). Agreeableness and Openness for White personas are essentially at the population average. In sum, aside from being more conscientious (and perhaps a touch more Machiavellian or anxious), White personas do not drastically differ from the mean persona profile in this dataset. Collectively, these profiles reinforce the 'model minority' myth for Asian characters—competent yet passive—while establishing White characters as the normative baseline with a capacity for strategic agency. This essentialist framing risks limiting narrative complexity, confining groups to predictable, culturally ingrained roles [11].

4.4. Gender-Driven stereotypes

Clear patterns of gender-based stereotyping emerge in the persona trait data. As we can see in Figure 7 (panel A) shows that female personas, on average, differ significantly from male personas on virtually every trait, with opposite-sign Z-scores for females vs. males in almost all cases. Female characters score higher on Agreeableness and Openness than their male counterparts, while scoring lower on Extraversion, Neuroticism, and all three Dark Triad traits. In numeric terms, the average female persona has A about +0.25 (in Z-score units) whereas the average male is around -0.25, and similarly O is about +0.3 for females versus -0.3 for males. This indicates the LLM often characterized women as more cooperative (high A) and more imaginative or open-minded (high O) than men. Conversely, female personas are portrayed as slightly more reserved on average (lower E) and—somewhat counterintuitively—far more emotionally stable (much lower N) than male personas. In fact, males in the dataset were depicted with a substantially higher Neuroticism (around +0.4) while females were around -0.4, meaning the model frequently made male characters more prone to stress or emotional volatility, whereas it cast female characters as unusually calm or emotionally steady. Conscientiousness is the one Big Five trait with only a slight gender difference: men were marginally above the mean and women marginally below, suggesting men were seen as just a bit more organized or disciplined, but this gap is very small.

All Dark Triad traits are strongly differentiated by gender in these personas. Men are assigned higher dark-trait scores across the board. On average, male personas score about 0.5–0.6 standard deviations higher in Machiavellianism than females (male M roughly +0.3 vs female M about -0.3). Likewise, male Psychopathy is

higher by roughly 0.36 z (male PSY around +0.18 vs female PSY -0.18). Narcissism shows a smaller gap (male NAR slightly above 0, female NAR slightly below 0), but even this difference is statistically reliable. These results indicate that the LLM frequently imbued male characters with more manipulative, self-focused, and callous traits compared to female characters, who were conversely depicted as less antagonistic and more pro-social.

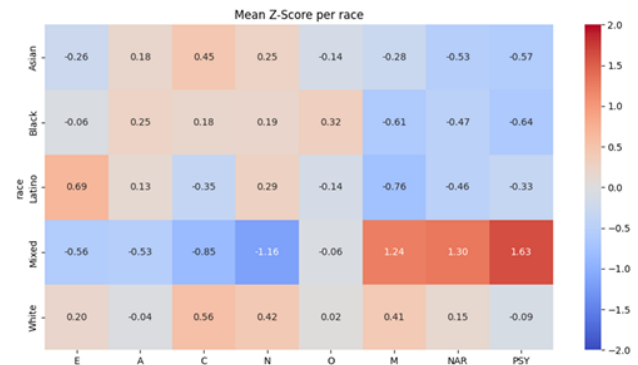


Figure 6: Mean Z-score per race. Heatmap of average standardized trait scores for personas of five racial categories (Asian, Black, Latino, Mixed, White). Trait abbreviations and color scale as before.

The visualization below in Figure 7 corroborates these differences. Panel B of Figure 7 displays the distribution of raw trait scores by gender, confirming systematic shifts: for each trait, the female distributions (orange boxplots) are centered at different levels than the male distributions (blue boxplots). For example, in Agreeableness, the female box is centered higher than the male box (most women personas scored more agreeable than most men), while in Neuroticism the male box is much higher than the female box (many male personas had high N scores, whereas female personas tended to have low N). Traits with large mean differences (like N, M, A) show clearly separated boxplot centers, whereas traits with smaller differences (like C, NAR) still have overlapping distributions but distinct averages. Panel C quantifies the mean gender differences (male minus female) in trait Z-scores with a bar chart. Each gray bar extending to the right indicates a higher male mean, and to the left a higher female mean; *p*-values from statistical tests are annotated. All traits show a significant difference ($p < 0.05$) between male and female personas. The largest gaps are observed in Neuroticism and Openness (males much higher in N, females much higher in O, both with $p < 0.001$), followed by Machiavellianism and Agreeableness (males higher in M, females in A, also highly significant). Psychopathy and Extraversion differences (males > females) are somewhat smaller but still clearly significant, and even the subtle differences in Conscientiousness and Narcissism reach significance. In sum, the persona dataset reveals a consistent gender-stereotypical pattern: male personas are generally portrayed as more extraverted, more neurotic, and higher on antagonistic/dark traits (M, NAR, PSY), whereas female personas are portrayed as more agreeable, more open, less neurotic, and lower on those dark traits.

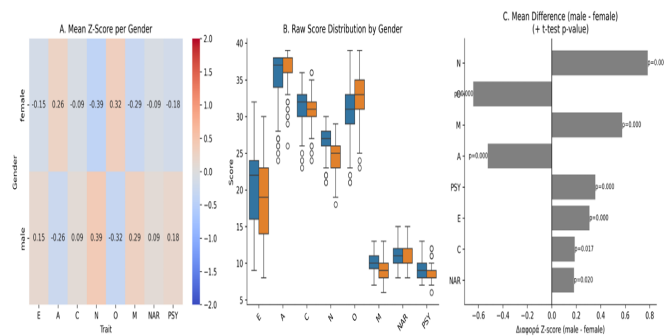


Figure 7: Gender differences in trait scores. Panel A – Heatmap of mean Z-scores for Female vs Male personas on each trait. Panel B – Boxplot distributions of raw trait scores by gender (blue = male, orange = female) for each trait (Big Five and Dark Triad). Panel C – Mean difference (male minus female) in Z-scores for each trait.

4.5. Intersections and PCA Clustering

To visualize how these trait biases combine and whether distinct demographic profiles cluster together, we performed a principal component analysis across all persona trait profiles. Figure 8 shows a scatter plot of all personas in the space of the first two principal components (PC1 vs PC2), with each point colored by race and marked by gender. Several clear patterns emerge. Race-based clustering is evident, particularly for the Mixed-race personas (purple points): they are widely separated from the rest, often occupying extreme positions in the plot. Many Mixed persona points lie far out on the rightmost end of PC1 or high on PC2, forming a distinct cloud largely isolated from other races. This reflects our earlier observation that Mixed-race profiles have extreme trait values (especially very high dark traits), which drive them to the periphery of the PCA space. For example, the cluster of purple symbols on the far right corresponds to Mixed personas with exceptionally high Machiavellianism/Narcissism/Psychopathy scores (traits likely loading heavily on PC1), while a subset of purple points that rise to the top of the chart represents Mixed personas that are outliers on a second combination of traits (perhaps those with unusual Big Five patterns contributing to a high PC2). A few of these extreme outliers are labeled by index in the figure, underscoring how far removed they are from the central mass of points.

In contrast, personas of other races (White, Black, Asian, Latino) tend to cluster nearer to the origin of the PCA plot and overlap considerably with each other. The dense central cloud of points (PC1 and PC2 values both near 0) is a mix of blue, orange, green, and red markers, indicating that White, Black, Asian, and Latino personas share a broadly similar trait space without forming wholly distinct clusters in the first two principal components. There are subtle tendencies—for instance, many Latino personas (red) appear slightly toward the left side of the central cluster (somewhat negative on PC1), whereas White (blue) and Asian (green) personas are more

dispersed around the middle, and Black personas (orange) intermingle throughout. However, these differences are gradual and overlapping; no single non-Mixed race forms an isolated grouping in this 2D projection. This suggests that aside from the Mixed category, racial trait differences are more a matter of degree than completely separate categories, with significant commonality among White, Black, Asian, and Latino personas in how the model represents their trait combinations.

Gender, indicated by shape (circles for male ● vs crosses for female ×), does not produce starkly separate clusters in the PCA plot. Male and female personas broadly overlap in this trait space, consistent with the fact that the gender differences we observed — although significant — involve opposing shifts on multiple traits that don't align neatly along a single principal axis. In Figure 6, male and female symbols of the same color are generally intermixed rather than split apart. For example, blue crosses and blue circles (female vs male White personas) are distributed in a similar area, and the same holds for other races (e.g., orange crosses and circles for Black personas largely coincide). This indicates that within each racial group, the gender-based trait offsets (e.g., females having slightly higher A and O, males higher M and N, etc.) add some scatter but do not create a separate “male persona cluster” distinct from a “female persona cluster.” The within-race variability — especially the extreme outlier status of certain races like Mixed — dominates the first two PCs.

That said, there are minor interaction effects visible. Within the Mixed-race cluster, female Mixed personas (purple ×) tend to concentrate a bit higher on the PC2 axis, whereas male Mixed personas (purple ●) extend further on PC1. This suggests that for Mixed-race characters, being male vs female leads to slightly different extreme trait manifestations: for instance, a Mixed male persona might combine the strong negative racial stereotype (Mixed: very low Big Five, very high dark traits) with the male-associated higher dark traits, yielding an especially extreme point far out on the PC1 dimension; a Mixed female, while still an outlier, may be somewhat tempered in dark traits (since females had lower dark scores) but could differ in another way (perhaps lower Neuroticism or higher emotional stability relative to Mixed males), pulling her profile in a slightly different direction (higher on PC2). Outside of the Mixed group, most other race-gender combinations do not produce clearly separable sub-clusters; the male-female differences within White, Black, Asian, and Latino groups appear as small shifts around a common central cluster for each race. Overall, the PCA visualization reinforces that race-based variations (the outlying nature of Mixed-race personas) are the primary driver of dispersion in trait space, while gender differences, though systematic, contribute more to fine-

scale variation within each racial cluster rather than forming entirely distinct groupings on the global map.

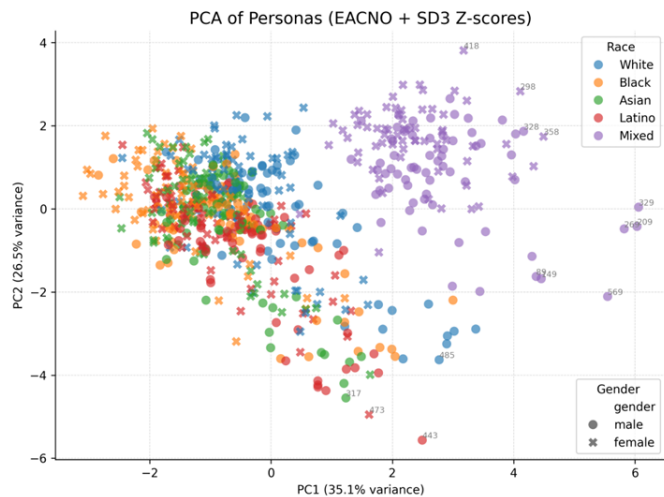


Figure 8: PCA of personas by race and gender. Scatter plot of persona trait profiles projected onto the first two principal components (PC1 and PC2, capturing ~61.6% of variance).

4.6. Internal Trait Correlations

The relationships among all the personality traits in this persona dataset provide insight into how traits tend to co-occur in the model's outputs. Figure 9 below shows the correlation matrix for every pair of traits. Several salient patterns stand out. Within the Big Five traits (the upper-left 5x5 block of the matrix), most correlations are positive, meaning that if a persona is high on one of these desirable traits, the model often also assigns higher levels on others. Notably, Agreeableness (A) strongly co-occurs with Openness (O) and Conscientiousness (C) (with Pearson r of roughly +0.70 for A–O and +0.54 for A–C). This indicates that more agreeable personas are also often portrayed as substantially more open-minded and responsible. Conscientiousness in turn has a moderate positive correlation with Openness ($r \approx +0.44$). These inter-correlations (A–C–O) suggest a “bundle” of positive traits in the dataset: many personas score high (or low) simultaneously on these three dimensions. Other Big Five pairs show weaker links; for example, Extraversion (E) is almost uncorrelated with Conscientiousness or Openness, and it has a slight negative correlation with Agreeableness (in this data, more extraverted characters were, if anything, a bit less agreeable, though the effect is small). Interestingly, Neuroticism (N) is nearly uncorrelated with most other Big Five traits here (its correlations with E, A, and C are close to zero). In short, aside from the cohesive cluster of A, C, and O moving together, the Big Five trait correlations are modest in magnitude.

By contrast, the Dark Triad traits show very strong mutual correlations. Machiavellianism, Narcissism, and Psychopathy are all positively interrelated, reflecting that personas who are high in one “dark” trait tend to be high in the others as well. The correlation between

Machiavellianism (M) and Psychopathy (PSY) is especially high ($r \approx +0.63$), and Machiavellianism also correlates around +0.60 with Narcissism (NAR). The NAR–PSY correlation is slightly lower (around +0.57) but still strong. This trio of high inter-correlations (the bright red block in the Dark Triad section of the matrix) indicates that the model often assigns all three dark traits in tandem — i.e. when it creates a manipulative persona, that character is also likely to be narcissistic and somewhat psychopathic in the portrayal. This is consistent with earlier observations that certain demographic groups (like Mixed-race or male personas) tended to receive uniformly high dark trait scores.

Looking at cross-domain relationships (Big Five vs. Dark Triad), we observe a clear inverse pattern between pro-social personality traits and the dark traits. Agreeableness has substantial negative correlations with Machiavellianism and Psychopathy ($r \approx -0.37$ and -0.41 , respectively). In other words, more agreeable (kind, empathetic) characters are much less likely to be portrayed as manipulative or callous. Conscientiousness likewise correlates negatively with Psychopathy (around -0.41), indicating that diligent, rule-abiding personas tend not to have psychopathic tendencies in the model's depiction. Neuroticism shows a moderately strong negative correlation with Narcissism ($r \approx -0.42$), suggesting that personas who are very narcissistic (self-important and confident) are often simultaneously depicted as emotionally stable (low N) rather than anxious — hinting that the model may associate narcissistic personalities with a kind of unshakeable confidence. Openness and Extraversion have weaker or mixed relationships with dark traits (most of those correlations hover near zero or a slight negative). One subtle finding is a slight positive correlation between Openness and Narcissism ($r \sim +0.17$), which implies that some highly open/intellectual personas were also given a hint of self-importance by the model. Additionally, Agreeableness versus Narcissism shows a very small positive r ($\sim +0.12$), meaning that unlike Machiavellianism and Psychopathy (which strongly conflict with Agreeableness), Narcissism in this dataset was not strongly anti-correlated with being agreeable — a persona could be somewhat agreeable and yet narcissistic (perhaps reflecting stereotypes of charming, sociable narcissists). Nonetheless, the dominant trend is that high dark-trait personas tend to score low on Agreeableness and Conscientiousness (seen in the blue-colored cells for A–M, A–PSY, C–PSY in Figure 7), reinforcing that benevolent personality characteristics are inversely related to antagonistic ones in the model's representation.

Overall, the correlation analysis confirms internally consistent patterns in the LLM's persona outputs. Positive personality traits align together and generally oppose the dark traits, while the Dark Triad traits form their own tight-knit cluster. These results provide a complementary

perspective on the trait structure underlying the demographic biases described above, demonstrating that the model's stereotypical persona attributions are not random but follow logical relationships (e.g., "kindness" versus "cruelty" as opposing poles, and certain positive traits tending to go hand-in-hand).

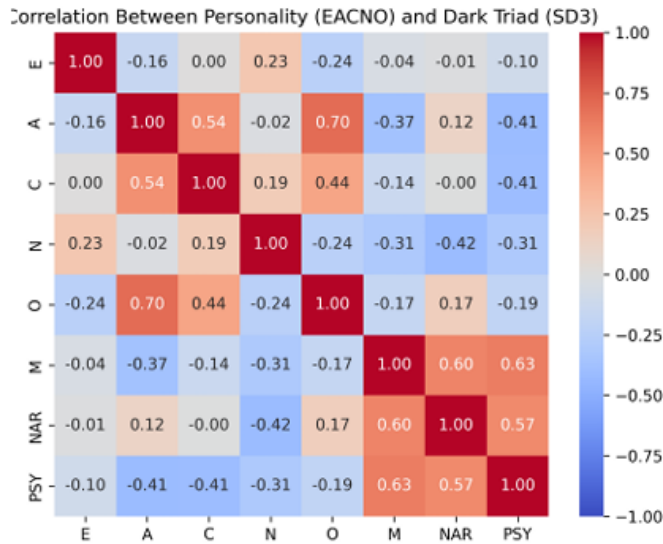


Figure 9: Correlation matrix of all traits. Pearson correlation coefficients between each pair of traits (Big Five: E, A, C, N, O; Dark Triad: M, NAR, PSY), computed across all persona scores. The matrix is symmetric; only one triangle is annotated with r values for clarity. Red indicates a positive correlation; blue indicates a negative correlation (scale shown on right).

5. Discussion

5.1. Cognitive and Psychological Interpretation

The observed patterns suggest that LLMs have developed internal cognitive-like representations of human groups, shaped by the statistical regularities of language. Although LLMs lack consciousness or intention, their training on vast human text corpora implicitly encodes societal narratives — producing what may be described as *synthetic cognition*. Unlike studies comparing AI to human baselines, our approach intentionally isolates this 'synthetic cognition' as a closed system. By focusing exclusively on the internal consistency of the model's generated personas, we map the algorithm's inherent stereotypical landscape without the confounding noise of human cultural variance.

The model's ability to assign coherent and demographically consistent personality profiles indicates that its latent representations capture more than linguistic associations: they embody social schemas. These schemas operate analogously to human stereotypes — simplifying complex social realities into categorical personality assumptions.

For instance:

- The "Western male atheist" archetype characterized by high *Openness* and *Narcissism*,

- The "Asian female Buddhist" with high *Conscientiousness* and low *Extraversion*, and
 - The "Black male Christian" with high *Extraversion* and *Agreeableness*
- demonstrate that the model generalizes culturally learned personality scripts.

Such patterns align with social cognition theory, which posits that stereotypes arise from heuristic associations rather than explicit reasoning. In this sense, the LLM functions as a large-scale mirror of human collective cognition — reproducing implicit personality prototypes learned from text [12].

5.2. Theoretical and Methodological Implications

From a methodological standpoint, this study bridges computational psychometrics and AI fairness auditing. Traditional bias research focuses on overt lexical or sentiment asymmetries (e.g., word embeddings associating "doctor" with male pronouns). Here, the bias operates at a *latent psychometric layer*, revealing how models attribute moral and emotional structure to demographic identities.

This framework contributes to the field by:

1. Introducing quantitative psychometric elicitation as a fairness diagnostic tool.
2. Demonstrating that *demographic conditioning* can alter inter-trait correlations — a deeper structural form of bias than mere mean-level differences.
3. Showing that bias can be interpreted through psychological theory, not just mathematical metrics.

Methodologically, it establishes a reproducible paradigm: using validated personality inventories (Big Five and Dark Triad), persona conditioning, and statistical normalization to extract interpretable cognitive maps from LLMs. This approach can be generalized to future studies exploring emotion, values, or moral reasoning biases in generative AI systems [13].

5.3. Ethical and Societal Considerations

The findings highlight serious ethical challenges. If LLMs systematically attribute moral or emotional traits based on identity cues, they risk reinforcing psychological stereotypes — subtle yet powerful forms of bias that influence downstream applications such as:

- Conversational AI: tone and empathy variation depending on user demographics;
- Hiring or profiling tools: skewed personality assessments;
- Education and therapy simulations: biased affective responses toward different identities.

- **Practical Applications of Psychometric Auditing:** Our framework could be extended to real-world applications beyond academic auditing. For example, it offers a method for monitoring racial bias trends in social media moderation systems, ensuring that automated agents do not attribute 'aggressive' or 'toxic' personality traits to users based on dialect or demographic markers. Furthermore, in the domain of healthcare, this methodology is critical for calibrating therapeutic LLMs. By detecting latent psychometric biases early, developers can fine-tune models to ensure they function equitably across diverse socio-economic and cultural backgrounds, preventing scenarios where an AI therapist might unconsciously adopt a colder or less empathetic persona toward marginalized groups."

Unlike explicit hate speech or toxicity, psychometric bias is invisible — it manifests through tone, moral emphasis, and perceived emotional intelligence. Because these models are often used in socially sensitive domains, their internal personality framing can affect fairness and trustworthiness.

To mitigate this, ethical AI development should include:

1. Psychometric fairness auditing — evaluating personality-related patterns alongside linguistic bias tests;
2. Data transparency — documenting sociocultural composition of training corpora;
3. Debiasing interventions — such as identity-neutral conditioning or fairness-aligned fine-tuning;
4. Human-in-the-loop oversight, ensuring that cultural interpretation does not reinforce stereotypes.

This work thus positions psychometric bias as a critical dimension of AI moral responsibility.

5.4. Limitations and Future Directions

Despite the robust methodology, several limitations must be acknowledged:

- **Synthetic Personas:** The personas simulate averaged demographic archetypes rather than real individuals, which limits ecological validity. However, this abstraction isolates model bias more effectively by removing user variance.
- **Single-Model Scope:** The experiments presented in the main analysis were conducted using one LLM (LLaMA-3.1-8B-Instruct). To assess whether the observed bias patterns are model-specific, we conducted preliminary exploratory experiments with additional models, including Mistral-7B-Instruct. These initial observations indicated qualitatively similar trends in demographic bias attribution, suggesting that the findings are not unique to a single

model architecture. However, a comprehensive cross-model validation, including proprietary models (e.g., GPT-4, Claude), is left as future work to determine the full extent of generalizability.

- **Cultural Bias in Training Data:** Because most pretraining text is in English, Western cultural norms dominate personality attributions. Extending this framework to multilingual LLMs could reveal cross-linguistic differences in psychometric stereotypes.
- **Simplified Gender Variable:** The binary male/female classification omits non-binary or gender-fluid identities, which may yield additional insight into model fairness.
- **Lack of Human Benchmark:** Although psychometric consistency was verified statistically, future work could compare LLM-generated profiles with human survey data to evaluate alignment.

Despite these limitations, the study establishes a foundational approach for examining how artificial cognition reflects human moral structure, offering a blueprint for next-generation bias auditing techniques [14], [6].

6. Conclusion and Future Work

This study introduced a novel framework for eliciting demographic stereotypes in Large Language Models (LLMs) through the lens of psychometric attribution. By combining established personality frameworks — the Big Five (EACNO) and the Dark Triad (SD3) — with systematic persona conditioning, we demonstrated that LLMs generate consistent, demographically structured personality profiles. These results provide compelling evidence that bias in LLMs extends beyond language or sentiment: it manifests at a cognitive level, where identity cues shape the model's perception of personality, morality, and social behavior.

Through large-scale experimentation across 660 personas, encompassing 11 regions, 5 racial groups, 6 religions, and 2 genders, the study revealed reproducible cross-group differences in both prosocial (Big Five) and antisocial (Dark Triad) traits. The model attributed:

- Higher *Agreeableness* and *Conscientiousness* to religious and female personas,
- Higher *Openness* and *Narcissism* to secular and Western personas,
- Greater *Machiavellianism* and *Emotional Restraint* to Asian personas,
- and elevated *Extraversion* and *Warmth* to African and Latin American personas.

These psychometric signatures were statistically significant and internally coherent, forming a structured "map of social cognition" embedded in the model's latent space.

In essence, the LLM acts as a mirror of collective cultural perception, reproducing personality stereotypes as learned from global human discourse.

From a theoretical standpoint, this work advances the field of computational psychometrics by framing model bias as a form of *synthetic cognition*. Rather than treating bias as a statistical defect, it reinterprets it as a *psychological phenomenon* — a window into how artificial systems internalize and reproduce the cognitive heuristics of human societies.

6.1. Key Contributions

1. **Methodological Innovation:** A reproducible Python-based pipeline for psychometric elicitation and statistical evaluation of demographic bias in LLMs.
2. **Theoretical Integration:** A bridge between AI fairness research, social psychology, and computational personality modeling.
3. **Empirical Findings:** Systematic personality and moral asymmetries across demographic factors, consistent with known cultural stereotypes.
4. **Ethical Insight:** Demonstration that fairness in LLMs must account for *psychological bias*, not only linguistic or representational bias.

6.2. Future Work

The present study opens several avenues for future research:

1. **Cross-Model Validation:** Extending the same pipeline to multiple LLM architectures (GPT-4, Claude, Gemini, Mistral) will reveal whether psychometric biases are *architecture-dependent* or *data-universal*.
2. **Temporal and Cultural Drift:** Investigating how model personality attributions evolve with new training data or fine-tuning cycles could expose *bias drift* over time.
3. **Multilingual and Cross-Lingual Evaluation:** Applying the framework to multilingual models may uncover differences in cultural stereotypes encoded across languages. This could lead to *comparative cultural cognition* analysis in AI.
4. **Inclusion of Non-Binary and Intersectional Identities:** Expanding demographic variables to include non-binary gender, mixed-religious backgrounds, and socioeconomic class will capture deeper intersectional complexity.
5. **Human Benchmarking:** Comparing LLM-generated profiles with actual psychometric data from human respondents can assess the degree of *alignment* between artificial and human stereotype structures.
6. **Bias Mitigation Techniques:** Implementing bias-aware fine-tuning, counter-stereotypical persona training,

and identity-neutral prompts could reduce psychometric distortion in model responses.

6.3. Final Remarks

The findings underscore a profound insight:

Large Language Models do not merely learn language — they learn society.

Their responses reveal a computational echo of human cognition, complete with virtues, flaws, and stereotypes. However, the implications of these findings reach far beyond technical correctness. As LLMs are increasingly integrated into decision-support systems for hiring, lending, and legal judgment, the implicit attribution of 'dark' or 'unstable' traits to specific demographics poses a tangible risk of algorithmic discrimination. If a model inherently views certain groups as less conscientious or more manipulative, this cognitive bias can cascade into material harm—denying opportunities or reinforcing systemic inequalities. Therefore, psychometric fairness is not merely a metric for model performance, but a safeguard for social justice in the age of artificial intelligence. The ultimate goal is to develop AI systems that reflect human diversity without reproducing human prejudice—systems that understand personality without imposing it. This study provides one step toward that vision, offering a reproducible foundation for exploring the psychology of artificial intelligence.

Ethical Disclosure

This research explicitly analyzes the generation of harmful stereotypes by AI systems. We acknowledge that some of the model-generated profiles reported—particularly those associating specific racial or religious groups with negative traits—contain offensive and discriminatory content. These outputs are presented solely for the purpose of scientific auditing and critique. The authors explicitly condemn these stereotypes and clarify that the demographic labels employed in this study (e.g., race, gender) are used as operational variables to probe the model's latent space, without implying essentialist definitions of complex human identities.

Conflict of Interest

The authors declare no conflict of interest.

Acknowledgment

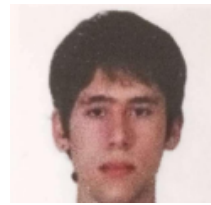
The author gratefully acknowledges the academic and technical support provided by colleagues and research collaborators during the design and implementation of this study. The experiments were conducted on locally maintained hardware resources, ensuring full reproducibility and data privacy.

No external funding was received for this work.

References

- [1] I. O. Gallegos et al., "Bias and fairness in large language models: A survey," *arXiv preprint arXiv:2309.00770*, 2023.
- [2] National Institute of Standards and Technology, "Towards a standard for identifying and managing bias in artificial intelligence," *NIST Special Publication 1270*, Gaithersburg, MD, 2023.
- [3] X. Bai, A. Wang, I. Sucholutsky, and T. L. Griffiths, "Explicitly unbiased large language models still form biased associations," *Proceedings of the National Academy of Sciences*, vol. 122, no. 8, p. e2416228122, 2025, doi: 10.1073/pnas.2416228122.
- [4] O. Gupta, S. Marrone, F. Gargiulo, R. Jaiswal, and L. Marassi, "Understanding social biases in large language models," *AI*, vol. 6, no. 5, p. 106, 2025, doi: 10.3390/ai6050106.
- [5] S. Lee et al., "Do LLMs have distinct and consistent personality? TRAIT: Personality testset designed for LLMs with psychometrics," in *Findings of the Association for Computational Linguistics: NAACL 2024*, 2024, doi: 10.48550/arXiv.2406.14703.
- [6] OpenAI, "GPT-4 Technical Report," *arXiv preprint arXiv:2303.08774*, 2023.
- [7] L. P. Argyle et al., "Out of one, many: Using language models to simulate human samples," *Political Analysis*, vol. 31, no. 3, pp. 337–351, 2023, doi: 10.1017/pan.2023.2.
- [8] D. Dodou, J. C. F. de Winter, and T. Driessen, "The use of ChatGPT for personality research: Administering questionnaires using generated personas," *Personality and Individual Differences*, vol. 228, p. 112729, 2024, doi: 10.1016/j.paid.2024.112729.
- [9] M. I. Radaideh, O. H. Kwon, and M. I. Radaideh, "Fairness and social bias quantification in large language models for sentiment analysis," *Knowledge-Based Systems*, vol. 319, p. 113569, 2025, doi: 10.1016/j.knosys.2025.113569.
- [10] D. S. Porat and E. Rabinovich, "Who are you, ChatGPT? Personality and demographic style in LLM-generated content," *arXiv preprint arXiv:2510.11434*, 2025.
- [11] S. Wang et al., "Exploring the impact of personality traits on LLM bias and toxicity," *arXiv preprint arXiv:2502.12566*, 2025.
- [12] H. Peters and S. C. Matz, "Large language models can infer psychological dispositions of social media users," *PNAS Nexus*, vol. 3, no. 6, p. pgae231, 2024, doi: 10.1093/pnasnexus/pgae231.
- [13] F. A. Tan et al., "PHAnToM: Persona-based prompting has an effect on theory-of-mind reasoning in large language models," in *Proceedings of the International AAAI Conference on Web and Social Media (ICWSM 2025)*, 2025.
- [14] T. Sühr, F. E. Dörner, S. Samadi, and A. Kelava, "Challenging the validity of personality tests for large language models," *arXiv preprint arXiv:2311.10805*, 2023.

Copyright: This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY-SA) license (<https://creativecommons.org/licenses/by-sa/4.0/>).



NIKOLAOS VASILEIOS OIKONOMOU is a Computer & Network Engineer, as well as an academic researcher and Ph.D. candidate in the Department of Informatics and Telecommunications at the University of Ioannina, from which he also received his B.Eng. and M.Sc. degrees. In parallel to his academic work, he serves as a private Computer Science educator and possesses several years of professional experience as a Software Developer, IT Specialist, and Network Consultant.



IOANNIS PALAIOKRASSAS is pursuing a M.Eng. degree in Computer Science and Engineering at the University of Ioannina and serves as an active research member. He is currently employed in web development.



DIMITRIOS VASILEIOS OIKONOMOU obtained his B.Sc. in Regional and Cross-Border Studies from the University of Western Macedonia in 2024. He is currently engaged in research activities at the same institution and is pursuing an M.Sc. in e-Business and Digital Marketing.



SOFIA PANAGIOTA CHALIASOU is pursuing a B.Sc. in Informatics at the Hellenic Open University and serves as an active research associate. She also holds a Vocational Diploma in Web Design and Development. In her professional capacity, she is currently employed in sales and possesses prior professional experience as a web developer.



NIKOLAOS RIGAS obtained his B.Sc. in Regional and Cross-Border Studies from the University of Western Macedonia in 2025. He is currently pursuing an M.Sc. in "Criminological and Penal Law perspectives on Corruption, Economic and Organized Crime" at the Hellenic Open University, while actively engaged in research activities.