# The Personality Trap: How LLMs Embed Bias When Generating Human-Like Personas

JACOPO AMIDEI, Universitat Oberta de Catalunya, Spain

GREGORIO FERREIRA, Universitat Oberta de Catalunya, Spain

MARIO MUÑOZ SERRANO, Universitat Oberta de Catalunya, Spain

RUBÉN NIETO, Universitat Oberta de Catalunya, Spain

ANDREAS KALTENBRUNNER, Universitat Pompeu Fabra, Spain

This paper examines biases in large language models (LLMs) when generating synthetic populations from responses to personality questionnaires. Using five LLMs, we first assess the representativeness and potential biases in the sociodemographic attributes of the generated personas, as well as their alignment with the intended personality traits. While LLMs successfully reproduce known correlations between personality and sociodemographic variables, all models exhibit pronounced WEIRD (western, educated, industrialized, rich and democratic) biases—favoring young, educated, white, heterosexual, Western individuals with centrist or progressive political views and secular or Christian beliefs. In a second analysis, we manipulate input traits to maximize Neuroticism and Psychoticism scores. Notably, when Psychoticism is maximized, several models produce an overrepresentation of non-binary and LGBTQ+ identities, raising concerns about stereotyping and the potential pathologization of marginalized groups. Our findings highlight both the potential and the risks of using LLMs to generate psychologically grounded synthetic populations.

Additional Key Words and Phrases: Bias, LLMs, Synthetic Population Sample, EPQR-A, Big Five personality inventory

## 1 Introduction

The use of large language models (LLMs) to generate synthetic populations offers a powerful, low-cost method for simulating human behavior. The ability to generate synthetic personas is being increasingly explored in fields ranging from psychology and political science to software engineering and healthcare. Synthetic personas, referred to by various names such as *guinea pigbots* [34], *silicon samples*, or *homo silicus* [33], have been proposed as substitutes for human participants in experimental research and as stand-ins for survey respondents or user testers. However, despite their growing popularity, important questions remain about how representative and unbiased these synthetic populations actually are. Recent evidence suggests that these representativeness concerns are not only a property of the underlying model, but also of how personas are elicited: persona prompting strategies can substantially change portrayals and stereotyping, and LLMs often struggle more when simulating marginalized groups [43]. Moreover, bias can be amplified when personas combine multiple traits in ways that are unlikely in real survey data ("incongruous personas"), leading to systematic deviations in persona-steered generation [40].

Synthetic populations are usually defined based on demographic [6, 20, 68], or individual-level input data [50] or giving information about historical or imaginary well-known persons [63], but here we investigate a different and under-explored method: generating populations based solely on personality test scores. The rationale for using personality traits stems from their well-established correlations with a wide range of social and behavioural outcomes, from moral judgments [5, 42, 61, 67] as well as value systems [15], civic engagement [66], preferences for and voting

for green parties [8], vaccination attitudes, intentions, and behaviours [9] to ethical vegetarianism and attitudes about animal welfare legislation [65, 70].

Leveraging these relationships, we examine how well LLMs can generate synthetic populations that reflect both the intended personality profiles and associated (if possible, unbiased) sociodemographic distributions. We pose three key research questions:

RQ1: Which sociodemographic attributes do LLMs generate when personas are based solely on personality test responses, and to what extent do these outputs reflect demographic biases?

RQ2: How does trait manipulation (e.g., maximizing Neuroticism or Psychoticism) influence demographic outputs?

RQ3: To what extent do LLM-generated personas internalize and express input personality traits?

To address these questions, we use responses from the Eysenck Personality Questionnaire Revised-Abbreviated (EPQR-A) [23] as inputs to five state-of-the-art LLMs and generate synthetic population samples.

To answer RQ1, we prompt the LLMs to provide sociodemographic attributes (e.g., gender, occupation, and political orientation) when generating the synthetic personas. Potential biases in the representativeness of the sample population can then be assessed by examining the distributions of these attributes. The same strategy was applied for RQ2, where the effect of extreme trait manipulation on the resulting generated personas is studied. Finally, to address RQ3, we employed multiple strategies: (1) measure the correlation between two comparable personality tests based on each synthetic persona's responses; (2) assess accuracy metrics by comparing synthetic responses to those of the input population; and (3) measure the internal consistency of the personality test employed.

Our findings suggest that LLMs are capable of generating synthetic sample populations that mirror to a large extent the input personality traits. However, the results also consistently reveal significant WEIRD biases. Synthetic populations skew heavily toward young, educated, white, Western, heterosexual individuals with secular or Christian backgrounds and progressive political views. More concerningly, we find that maximizing traits like Psychoticism leads to a marked overrepresentation of LGBTQ+ and non-binary identities in several models, raising important questions about harmful stereotypes and potential pathologization.

By combining psychometric rigor with generative analysis, this study contributes new insights into both the capabilities and risks of using LLMs to simulate human populations, especially when these models are treated as proxies for real-world diversity.

## 2  Related work

This paper bridges three research areas in the LLM framework: i) Building synthetic sample populations with LLMs, ii) Assessing LLM personality through the use of personality tests, and iii) Bias detection in LLMs.

### 2.1  Building synthetic sample populations with LLMs

The employment of LLMs as substitutes for human participants was studied in psychological research [17, 34, 52], political polling [58], software engineering research [28], teaching research [45], economics [33], social media platforms design [50, 69], market research to understand consumer preferences [10] and more generally social science research [6]. Across these domains, findings are mixed: some studies report alignment with human response patterns [6, 33, 44, 50, 51, 58, 68], whereas others caution that LLM substitutes can fail in systematic ways and should not be treated as drop-in replacements for human participants [2, 3, 13, 32, 52, 55, 71]. For example, Agnew et al. [2] sheds light on potential obstacles when utilizing LLMs to simulate human behaviour, such as the current inability of LLMs to emulate human cognition and

decision-making accurately, the reliance of psychology research on various non-linguistic cues to study human cognition and behaviour, and the phenomenon of "value lock-in" (that is the LLMs ability of reflecting attitudes only from the time of their training).

Recent work has also begun to test whether synthetic participants generalize beyond WEIRD settings and policy domains [64], and methodological discussions have started to formalize LLMs as artificial research participants and clarify the risks of substituting humans with model-generated respondents in behavioral research designs [47].

## 2.2 Assessing LLM personality through the use of personality tests

In recent years, there has been a surge in the employment of personality questionnaires within the LLM framework. For example, the Big Five factors [16] were used, among others, by Jiang et al. [35], Karra et al. [37], Mei et al. [48], Pellert et al. [54], Serapio-García et al. [62] to quantify the personality traits of LLMs. Similarly, IPIP-NEO [29] was used in Serapio-García et al. [62] and Short Dark Tetrad (SD4) [53] was used in Pellert et al. [54]. The EPQR-A was instead used in a multilingual setting in Amidei et al. [4], Ferreira et al. [20, 22]. While the outcomes of the aforementioned studies may vary depending on the LLMs and questionnaires used, there is support to conclude that personality assessments for LLMs are valid and reliable. Nevertheless, Dorner et al. [18], Gupta et al. [31], Zou et al. [73] argue against the use of self-reported tests for LLMs.

An important validity concern in questionnaire-based profiling is that LLMs may treat survey administration as an evaluation context and respond in systematically self-presentational ways. Salecha et al. [57] show that across multiple LLM families, responses to Big Five surveys shift toward socially desirable trait poles once models can infer that they are being evaluated. This line of work suggests that part of what psychometric instruments measure in LLMs can reflect emergent social-desirability responding rather than stable underlying traits, which is particularly salient for instruments (such as EPQR-A) that include explicit social-desirability components.

## 2.3 Bias detection in LLMs

As LLMs are increasingly used to generate text that stands in for people (e.g., synthetic respondents or personas), auditing demographic bias and representational harms becomes essential. In this work, we focus specifically on bias that emerges in LLM-generated synthetic populations. For a more comprehensive discussion on bias in LLMs, we refer the reader to [25, 30, 39]. Concerns about the use of synthetic sample populations, due to inherent biases in LLMs, have been raised by, among others, in [2, 13, 32, 71]. For example, Crockett and Messeri [13] points up the problem of population generability. This problem was studied by Harding et al. [32], who criticized the fact that LLMs' population representativeness must be carefully circumscribed. Wang et al. [71] further show that LLMs misportray and flatten demographic groups due to intrinsic model bias and limitations, including their difficulty in representing minority groups [2]. Indeed, the problem of misrepresentation, defined as an incomplete or non-representative distribution of a sample population generalised to a broader social group, is a common form of social bias identified in NLP [25].

Relatedly, recent NLP work has studied bias under persona-steered generation, showing that bias is not only a property of model parameters but also of how personas are composed and prompted. Liu et al. [40] introduce the notion of incongruous personas (multi-trait personas whose traits are unlikely to co-occur in human data) and show that LLM generations can deviate from expected distributions when personas combine multiple identity-relevant attributes. Complementing this, Lutz et al. [43] systematically evaluate sociodemographic persona prompting strategies and find that seemingly minor prompt choices (e.g., role-adoption format or demographic priming) can substantially change portrayals and stereotyping, particularly for marginalized groups. Finally, Ostrow and Lopez [49] document that

LLMs reproduce stereotypes about sexual and gender minorities beyond binary gender categories, both in survey-style elicitation and in free-form text generation, raising concerns about representational harms.

## 2.4 Research gap

To the best of our knowledge, prior work has not examined the generation of synthetic populations conditioned solely on responses (or scores) from personality questionnaires. This gap matters because personality-test conditioning is an appealingly lightweight way to generate large-scale personas without collecting biographical data. Yet it may also create a direct channel for stereotype activation when downstream demographic attributes are inferred from the personality-conditioned text. We therefore integrate the three strands above to evaluate both (i) whether personality-score conditioning yields coherent and diverse personas, and (ii) whether it introduces systematic demographic distortions and representational harms.

## 3 Methods

### 3.1 Experimental Setups

We designed a multi-step experimental pipeline combining psychometric profiling, generative prompting, and statistical evaluation.[1] The different steps of our pipeline are visualized in Figure 1.

The input to our study is a dataset of 826 simulated responses to the EPQR-A personality questionnaire, provided by Ferreira et al. [21][2]. The authors of this study used different questionnaire languages and parameter settings. Here we only use the responses to the English version of the EPQR-A completed by gpt-4o-2024-05-13 with temperature 1. Each of these responses was used as input to prompt an LLM to generate a detailed synthetic persona, which should represent a synthetic persona whose demographic, biographical, and behavioral features were to reflect the personality traits inferred from the questionnaire.

We tested five large language models: GPT-3.5 (gpt3.5-turbo-0125), GPT-4o (gpt4o-2024-11-20), claude-3.5-s (claude3.5-sonnet), and the two open-source models LLaMa3.2-3B and LLaMa3.1-70B.[3] Each model was prompted with a consistent template that embedded the EPQR-A responses and instructed the model to imagine and describe a person based on that personality profile. The output was required in a structured JSON format containing a predefined set of 8 sociodemographic attributes: age, gender, sexual orientation, race, religious belief, occupation, political orientation, and location. See Appendix B.1 for the corresponding prompts and Appendix A for examples of the resulting synthetic persona descriptions.

In addition to this unaltered or "baseline" generation condition (to which we refer to as Base population in the remainder of this study), we generated a random control baseline condition and two manipulated conditions (MaxN and MaxP) to assess how models respond to extreme personality trait anchoring. The random condition was generated by computing the true/false binomial distributions for each question and then randomly generating 826 answered questionnaires using these distributions. For the MaxN condition, we adjusted each input response to yield the highest possible Neuroticism score while keeping other traits unchanged; similarly, for the MaxP condition, we maximized Psychoticism. MaxN and MaxP enabled us to investigate how exaggerated personality traits impact sociodemographic outcomes and whether they elicit stereotypical or pathologising patterns.

---

[1]The code used for our experiments can be found at https://anonymous.4open.science/r/the_personality_trap-F487/README.md.

[2]The synthetic populations generated in Ferreira et al. [21] were based on the demographic characteristics (655 females: mean age = 18.9, SD = 1.56 and 171 males: mean age = 19.36, SD = 1.99) of a population described in García-González et al. [26].

[3]We performed our experiments for the GPT model family with OpenAI's API (total cost of 228.12 USD) and for the remaining models on Amazon AWS (total cost of 223.96 USD).
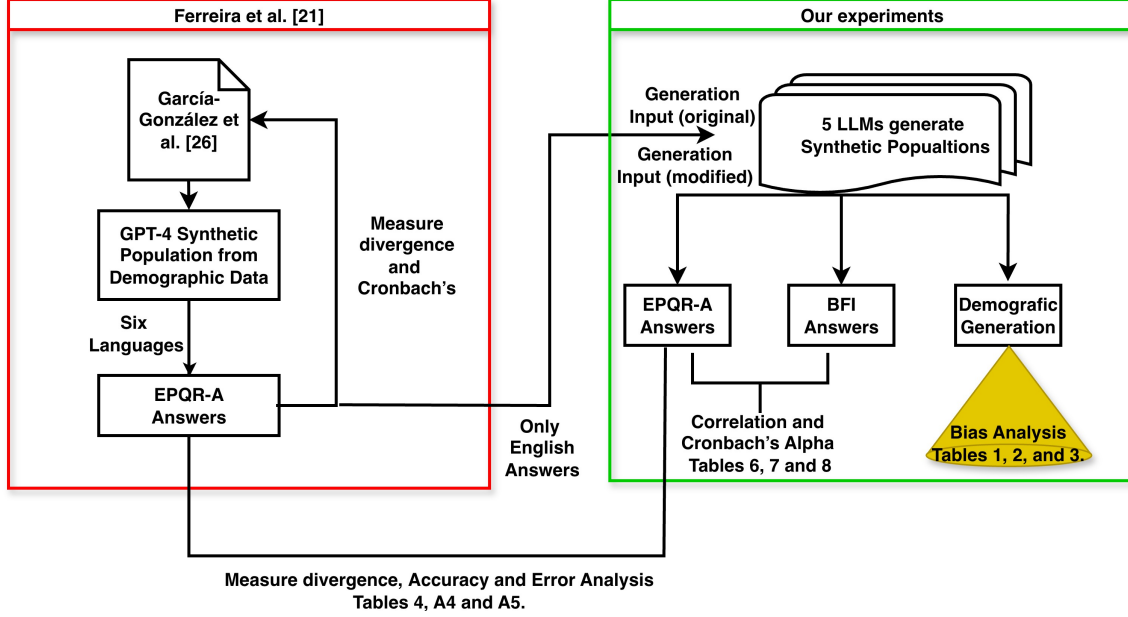
Fig. 1. Overview of the experimental pipeline. Our experiment (right square) builds on the English version of the EPQR-A responses from 826 simulated personas reported in Ferreira et al. [21] (left square), which were used as generation input. Five LLMs were then prompted to generate 826 synthetic personas and provide sociodemographic attributes while reflecting the input personality traits. Subsequently, the LLMs were tasked to answer both the EPQR-A and BFI with the generated personalities. In a further step, the LLMs generated another 826 synthetic personas based on extreme trait manipulations, again producing sociodemographic attributes consistent with the input traits. Potential biases in the representativeness of the sample populations were assessed by analysing the distributions of these attributes. Divergence, accuracy, and error analysis were computed between the input responses and those newly generated. In addition, correlations between the EPQR-A and the BFI were examined, along with internal consistency measured using Cronbach's alpha.

## 3.2 Strategies for Result Evaluation

To evaluate generation consistency, we generated 10 Base sample populations and 5 MaxN and MaxP sample populations using identical personality measures as inputs (i.e., answers to the EPQR-A questionnaire). To determine whether the generated personas reflected meaningful and coherent demographic patterns (RQ1 and RQ2), we analyzed five key identity categories (Gender, Race, Religion, Political, and Sexual Orientation) using descriptive statistics and two-sided t-tests to compare means between Base and MaxN and MaxP populations.

To assess trait fidelity and alignment (RQ3), we conducted another round of testing in which the generated personas were asked to complete the EPQR-A questionnaire again. Since the Inter-sample variation was minimal, we selected one representative population sample from each condition (Base, MaxN, and MaxP) for this analysis. We then computed accuracy metrics such as Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) between the input and regenerated EPQR-A scores. For the Base sample, we also administered the Big Five personality inventory (BFI) [36]. Pearson correlations between EPQR-A and BFI scores were calculated to examine cross-instrument consistency, and Cronbach's alpha [14] was used to assess internal reliability for both EPQR-A and the BFI.

Table 1. Averages and standard deviations of ten Base, five MaxN, and five MaxP sample populations statistics for GPT-4o and GPT-3.5. Non-bin., Con., Prog. Hetero. and Unspe. stands for non-binary, conservative, progressive, heterosexual, and unspecified, respectively. Cases where results are significantly different from the Base sample populations are marked as p<0.05*, p <0.01†, and p <0.001‡ (two-sided t-test).

| | | GPT-4o | | | GPT-3.5 | | |
|---|---|---|---|---|---|---|---|
| | | Base | MaxN | MaxP | Base | MaxN | MaxP |
| Gender | Female | 25.71 ± 1.36 | 29.54 ± 2.07‡ | 4.67 ± 0.79‡ | 90.82 ± 2.13 | 90.94 ± 1.21 | 88.74 ± 1.02† |
| | Male | 44.75 ± 2.42 | 30.46 ± 2.33‡ | 6.27 ± 0.89‡ | 9.03 ± 2.28 | 8.86 ± 1.14 | 10.94 ± 1.00† |
| | Non-bin. | 29.18 ± 1.76 | 39.71 ± 1.27‡ | 88.76 ± 0.99‡ | 0.14 ± 0.21 | 0.19 ± 0.14 | 0.31 ± 0.18 |
| | Other | 0.36 ± 0.21 | 0.29 ± 0.22 | 0.29 ± 0.20 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 |
| Pol. Or. | Centre | 64.50 ± 1.43 | 51.45 ± 1.47‡ | 0.36 ± 0.17‡ | 75.08 ± 0.89 | 72.98 ± 1.68* | 18.23 ± 1.13‡ |
| | Con. | 0.12 ± 0.08 | 0.02 ± 0.05 | 0.00 ± 0.00* | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 |
| | Prog. | 34.63 ± 1.44 | 48.11 ± 1.46‡ | 97.70 ± 0.19‡ | 16.46 ± 1.79 | 18.55 ± 1.56‡ | 38.74 ± 1.95‡ |
| | Others | 0.75 ± 0.13 | 0.41 ± 0.22* | 1.94 ± 0.21‡ | 8.45 ± 1.68 | 8.48 ± 0.67 | 43.03 ± 1.39‡ |
| Race | Asian | 0.58 ± 0.29 | 0.63 ± 0.34 | 1.89 ± 0.53‡ | 3.10 ± 0.41 | 2.88 ± 0.40 | 2.30 ± 0.24* |
| | Black | 0.05 ± 0.07 | 0.00 ± 0.00 | 0.29 ± 0.18† | 0.97 ± 0.33 | 0.87 ± 0.16 | 1.16 ± 0.55 |
| | Latin | 0.12 ± 0.00 | 0.19 ± 0.11 | 0.46 ± 0.23† | 7.33 ± 1.64 | 8.26 ± 0.96 | 6.10 ± 0.66* |
| | White | 98.98 ± 0.43 | 98.64 ± 0.43 | 90.39 ± 1.00‡ | 88.59 ± 1.83 | 87.94 ± 0.83 | 90.41 ± 0.58† |
| | Other | 0.26 ± 0.18 | 0.53 ± 0.14 | 6.98 ± 0.77‡ | 0.00 ± 0.00 | 0.05 ± 0.11 | 0.02 ± 0.05 |
| Relig. Beli. | Christian | 12.86 ± 1.35 | 11.62 ± 0.65 | 0.00 ± 0.00‡ | 3.83 ± 0.90 | 4.72 ± 0.86* | 0.46 ± 0.05‡ |
| | Agnostic | 84.74 ± 1.16 | 84.29 ± 0.55 | 94.60 ± 0.41‡ | 91.06 ± 0.72 | 90.68 ± 0.89 | 93.63 ± 0.64‡ |
| | Atheist | 0.12 ± 0.12 | 0.00 ± 0.00* | 4.02 ± 0.47‡ | 3.66 ± 1.05 | 2.90 ± 0.42 | 5.25 ± 0.73‡ |
| | Others | 2.28 ± 0.48 | 4.09 ± 0.26‡ | 1.38 ± 0.27† | 1.45 ± 0.27 | 1.69 ± 0.31 | 0.66 ± 0.20‡ |
| Sex. Or. | Hetero. | 70.07 ± 1.63 | 59.57 ± 1.02‡ | 5.88 ± 0.79‡ | 99.86 ± 0.21 | 99.81 ± 0.14 | 99.66 ± 0.16 |
| | LGBTQ+ | 29.93 ± 1.63 | 40.41 ± 1.07‡ | 94.12 ± 0.79‡ | 0.14 ± 0.21 | 0.19 ± 0.14 | 0.31 ± 0.18 |
| | Unspe. | 0.00 ± 0.00 | 0.02 ± 0.05 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.02 ± 0.05 |

### 3.3 Personality tests used in the experiments

The EPQR-A is an abbreviated version of the Eysenck Personality Inventory [19], containing 24 items to assess three personality dimensions (6 items each): **Extraversion (E), Neuroticism (N), and Psychoticism (P)**. It also includes a scale to assess social desirability **Lie (L)**, which also contains six items. Each item has a dichotomous response (yes or no), and a score for each scale can be computed by summing individual items (resulting in a range from 0 to 6).

The BFI contains 44 items (assessed in a scale from 1 to 5) that measure individuals in the following general dimensions: **Extraversion (E), Agreeableness (A), Conscientiousness (C), Neuroticism (N)**, and **Openness (O)**. Both the EPQR-A and the BFI assess N and E, and taking into account available results in the literature, we hypothesize that scores from the corresponding dimensions of the two questionnaires should be highly and significantly correlated [7]. Conversely, the dimensions C and A of the BFI are expected to be inversely correlated with the dimension P from the EPQR-A, while the dimension O from the BFI should be positively correlated with E from the EPQR-A [56].

### 4 Results

### 4.1 Analysing sociodemographic attributes generation (RQ1 and RQ2)

We first ask what sociodemographic characteristics LLMs assume when asked to generate personas from personality-test responses alone (RQ1), and how these assumptions shift under extreme trait manipulation (RQ2).

Tables 1, 2, and 3 report the mean and standard deviation (calculated across ten trials for the Base populations and five trials each for the MaxP and MaxN populations) across five sociodemographic categories: Gender, Political Orientation,

Table 2. Averages and standard deviations of ten Base, five MaxN and five MaxP sample populations statistics for LLaMa3.2-3B and LLaMa3.1-70B. Cases where the results are significantly different from the Base sample populations are marked as p<0.05*, p <0.01†, and p <0.001‡ (two-sided t-test). Abbreviations as in Table 1.

| | | LLaMa3.2-3 | | | LLaMa3.1-70B | | |
|---|---|---|---|---|---|---|---|
| | | Base | MaxN | MaxP | Base | MaxN | MaxP |
| Gender | Female | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 22.03 ± 2.56 | 30.07 ± 0.40‡ | 5.04 ± 0.90‡ |
| | Male | 100.00 ± 0.00 | 100.00 ± 0.00 | 100.00 ± 0.00 | 77.94 ± 2.56 | 69.81 ± 0.28‡ | 94.34 ± 0.85‡ |
| | Non-bin. | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.02 ± 0.05 | 0.02 ± 0.05 | 0.38 ± 0.16‡ |
| | Other | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.10 ± 0.10* | 0.24 ± 0.23† |
| Pol. Or. | Centre | 1.94 ± 0.33 | 2.16 ± 0.68 | 0.00 ± 0.00‡ | 10.41 ± 0.64 | 6.20 ± 0.61‡ | 0.00 ± 0.00‡ |
| | Con. | 32.56 ± 1.68 | 15.35 ± 1.03‡ | 0.00 ± 0.00‡ | 42.93 ± 1.23 | 42.16 ± 0.75 | 0.00 ± 0.00‡ |
| | Prog. | 65.37 ± 1.57 | 82.42 ± 1.13‡ | 98.67 ± 0.52‡ | 37.19 ± 1.69 | 50.43 ± 0.67‡ | 69.88 ± 1.23‡ |
| | Others | 0.12 ± 0.12 | 0.07 ± 0.11 | 1.33 ± 0.52‡ | 9.47 ± 0.56 | 1.21 ± 0.35‡ | 30.12 ± 1.23‡ |
| Race | Asian | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.90 ± 0.80 | 0.36 ± 0.19† | 0.00 ± 0.00‡ |
| | Black | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 |
| | Latin | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.02 ± 0.05 | 0.00 ± 0.00 |
| | White | 100.00 ± 0.00 | 100.00 ± 0.00 | 100.00 ± 0.00 | 99.08 ± 0.79 | 99.52 ± 0.09* | 99.39 ± 0.31 |
| | Other | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.02 ± 0.05 | 0.10 ± 0.10 | 0.61 ± 0.31‡ |
| Relig. Beli. | Christian | 76.27 ± 0.78 | 72.64 ± 1.94‡ | 11.45 ± 1.63‡ | 61.16 ± 0.45 | 61.72 ± 1.35 | 0.00 ± 0.00‡ |
| | Agnostic | 23.24 ± 0.82 | 26.51 ± 1.79‡ | 67.75 ± 2.40‡ | 36.85 ± 0.90 | 37.38 ± 1.18 | 20.15 ± 2.03‡ |
| | Atheist | 0.48 ± 0.19 | 0.85 ± 0.27* | 20.80 ± 1.23‡ | 1.89 ± 1.08 | 0.87 ± 0.31‡ | 79.85 ± 2.03‡ |
| | Others | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.10 ± 0.10 | 0.02 ± 0.05 | 0.00 ± 0.00* |
| Sex. Or. | Hetero. | 100.00 ± 0.00 | 100.00 ± 0.00 | 100.00 ± 0.00 | 99.98 ± 0.05 | 99.90 ± 0.10 | 96.76 ± 0.52‡ |
| | LGBTQ+ | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.02 ± 0.05 | 0.00 ± 0.00 | 3.00 ± 0.44‡ |
| | Unspe. | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.10 ± 0.10* | 0.24 ± 0.23† |

Race, Religious Belief, and Sexual Orientation.[4] Three additional categories, that is, age, location, and occupation, were measured, but due to space constraints and since differences across populations (Base, MaxN, MaxP) and models were minimal, detailed results are not shown. However, relevant patterns are discussed throughout the paper.

We began by assessing generation consistency, measured as the variance across trials. The low standard deviations observed across categories in Tables 1, 2, and 3 indicate a high degree of consistency for each model.

*Analysing the Base synthetic populations (RQ1):* Examining Tables 1, 2, and 3 reveals that the synthetic sample populations generated by the models share several key traits: most are either male or female, with claude-3.5-s and GPT-3.5 skewing female, LLaMa3.2-3B and LLaMa3.1-70B models skewing male, and GPT-4o showing the most gender diversity, including 29.27% (± 1.70) non-binary individuals. Racially, the majority are white except for claude-3.5-s, which includes 35.16% (± 1.72) Asian. Most are heterosexual, except for GPT-4o, which includes 29.96% (± 1.52) LGBTQ+. Regarding religious beliefs, all models (but to a lesser extent for the LLaMa family, which skew towards Christians), tend to be agnostic or outside major faiths like Buddhism and Hinduism, with no instances of Islam and Judaism. Politically, they lean centrist or progressive, though LLaMa3.2-3B and LLaMa3.1-70B models include significant conservative portions, that is 42.68% (± 1.00) and 33.16% (± 1.30) respectively. Furthermore, the analysis of age, location, and occupation reveals that the Base population primarily consists of individuals aged on average 28 to 32, based mostly in major U.S. metropolitan areas (with London (UK) as the only exception), such as Chicago, New York City, San Francisco, and Los Angeles. They are typically employed in highly educated fields, including Accounting & Finance, Tech & Engineering, Creative & Design, Research & Science, and Health & Social Care.

---

[4]To analyze model outputs consistently across experiments/trials and models, we normalized the free-text sociodemographic attributes produced by the LLMs into a compact, pre-defined set of categories. More details in Appendix C.

Table 3. Averages and standard deviations of ten Base, five MaxN, and five MaxP sample populations statistics for claude-3.5-s. Cases where the results are significantly different from the Base sample populations are marked as p<0.05*, p <0.01†, and p <0.001‡ (two-sided t-test). Abbreviations as in Table 1.

| | | claude-3.5-s | | |
|---|---|---|---|---|
| | | Base | MaxN | MaxP |
| **Gender** | Female | 87.84 ± 1.09 | 98.31 ± 0.18‡ | 0.83 ± 0.10‡ |
| | Male | 8.93 ± 1.19 | 0.29 ± 0.11‡ | 2.52 ± 0.45‡ |
| | Non-bin. | 3.22 ± 0.18 | 1.40 ± 0.18‡ | 96.66 ± 0.48‡ |
| | Other | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 |
| **Pol. Or.** | Centre | 90.70 ± 1.23 | 79.18 ± 1.51‡ | 0.00 ± 0.00‡ |
| | Con. | 0.56 ± 0.21 | 0.00 ± 0.00‡ | 0.00 ± 0.00‡ |
| | Prog. | 8.67 ± 1.07 | 20.82 ± 1.51‡ | 75.72 ± 0.38‡ |
| | Others | 0.07 ± 0.07 | 0.00 ± 0.00 | 24.28 ± 0.38‡ |
| **Race** | Asian | 34.82 ± 1.43 | 29.15 ± 1.58‡ | 11.72 ± 0.42‡ |
| | Black | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 |
| | Latin | 0.12 ± 0.08 | 0.22 ± 0.16 | 0.05 ± 0.07 |
| | White | 62.03 ± 1.70 | 68.96 ± 1.79‡ | 6.13 ± 0.43‡ |
| | Other | 3.03 ± 0.72 | 1.67 ± 0.22‡ | 82.10 ± 0.63‡ |
| **Relig. Beli.** | Christian | 0.07 ± 0.11 | 0.05 ± 0.07 | 0.00 ± 0.00 |
| | Agnostic | 99.90 ± 0.10 | 99.95 ± 0.07 | 99.81 ± 0.16 |
| | Atheist | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.19 ± 0.16† |
| | Others | 0.02 ± 0.05 | 0.00 ± 0.00 | 0.00 ± 0.00 |
| **Sex. Or.** | Hetero. | 96.78 ± 0.18 | 98.60 ± 0.18‡ | 0.92 ± 0.14‡ |
| | LGBTQ+ | 3.22 ± 0.18 | 1.40 ± 0.18‡ | 99.08 ± 0.14‡ |
| | Unspe. | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.00 ± 0.00 |

*Comparing Base, MaxN and MaxP (RQ2):* We now analyse the synthetic populations generated when maximising the P or N scale of the answers to the input questionnaires.

When comparing the Base, MaxN, and MaxP synthetic sample populations, significant (based on a two-sided t-test) sociodemographic shifts emerge. In detail, the MaxN sample populations show more moderate shifts compared to MaxP, but still share some common trends, notably a general increase in progressive political alignment and a modest increase in occupations related to creativity or education. Claude-3.5-s slightly increases female representation to 98.31% (± 0.17), approximately doubles progressive presence, and slightly shifts from Tech-related professions to Education-related ones, while GPT-4o shows a milder gender shift with small increases in female and non-binary identities, and a rise in writing & publishing related professional occupations. LLaMa3.2-3B and LLaMa3.1-70B both boost progressive alignment (by ≈18% and ≈13%, respectively). GPT-3.5 is the most conservative in MaxN, showing minimal deviation from the Base sample population, making it an outlier in terms of sociodemographic change.

MaxP outputs consistently exhibit a marked shift toward progressive political alignment and increased representation of creative occupations, particularly in creative & design and writing & publishing roles. Notably, GPT-4o and Claude-3.5-S increase the prevalence of LGBTQ+ and non-binary identities; while this does not imply any real-world association, it raises concerns about stereotype-driven generation and potential pathologizing inferences. LLaMa3.1-70B, in contrast, shifts toward a more male-dominated distribution. Geographically, MaxP samples across models tend to cluster personas in Western and Southern U.S. cities. GPT-3.5 stands out for reducing creative & design roles and emphasising events & community-related professions instead; LLaMa3.2-3B and LLaMa3.1-70B models are distinct in driving strong religious shifts in detriment of Christianity, LLaMa3.2-3B toward agnosticism (similarly, but to a lesser extent, GPT-4o), and LLaMa3.1-70B toward Atheism.

Table 4. EPQR-A scores per model, sample population, and category. Significant differences from the input scores at the individual level: p<0.05* and p <0.01† (two-sided paired t-test) and at the population level: p<0.05§ and p<0.01¶ (two-sided unpaired t-test).

| Model | Pop. | EPQR-A | | | |
|---|---|---|---|---|---|
| | | E | N | P | L |
| GPT-4o | Base | 2.16±2.85† | 3.06±2.59 | 0.87±1.04 | 5.92±0.48† |
| | MaxN | 2.15±2.86† | 5.96±0.29† | 0.91±1.11† | 5.91±0.57* |
| | MaxP | 2.17±2.87† | 3.01±2.73 | 4.75±1.03† | 5.85±0.72 |
| GPT-3.5 | Base | 2.52±2.64†, § | 3.32±2.36†, § | 0.96±0.93†, § | 5.93±0.41† |
| | MaxN | 2.94±2.62†, ¶ | 4.05±1.76†, ¶ | 1.17±0.95†, ¶ | 5.94±0.38†, § |
| | MaxP | 3.71±2.32†, ¶ | 2.61±2.30†, ¶ | 3.12±1.18†, ¶ | 5.88±0.49 |
| claude-3.5-s | Base | 2.19±2.85† | 3.31±2.44† | 0.61±0.91†, ¶ | 5.90±0.60 |
| | MaxN | 2.18±2.86† | 6.00±0.00†, ¶ | 0.72±0.97†, ¶ | 5.87±0.61 |
| | MaxP | 2.22±2.85* | 3.15±2.68* | 5.51±0.65†, ¶ | 5.42±1.42†, ¶ |
| LLaMa3.2-3B | Base | 2.55±2.48†, § | 2.42±2.18†, ¶ | 1.38±1.11†, ¶ | 5.22±0.44†, ¶ |
| | MaxN | 2.48±2.55† | 3.94±2.06†, ¶ | 1.49±1.05†, ¶ | 5.04±0.29†, ¶ |
| | MaxP | 3.06±2.36†, ¶ | 2.72±2.22†, ¶ | 2.75±0.97†, ¶ | 4.99±0.31†, ¶ |
| LLaMa3.1-70B | Base | 2.22±2.83* | 3.49±2.33†, ¶ | 0.75±1.04†, § | 5.94±0.57† |
| | MaxN | 2.24±2.85 | 5.91±0.33†, ¶ | 0.75±0.97†, § | 5.95±0.53†, § |
| | MaxP | 2.54±2.77†, § | 2.99±2.60* | 4.50±1.46†, ¶ | 5.90±0.59 |
| | Input | 2.26±2.79 | 3.08±2.42 | 0.85±0.97 | 5.89±0.54 |
| | Random | 2.23±1.20 | 3.01±1.13 | 0.82±0.72 | 5.89±0.33 |

## 4.2 Analysing the degree to which LLM-generated personas reflect the intended personality traits (RQ3)

Despite large shifts in inferred demographics under trait manipulation, the personas generally preserve the intended personality profiles when evaluated through the same questionnaire that generated them. We test this trait fidelity by asking each persona to complete the EPQR-A again and comparing regenerated scores to the input distribution (MAE/RMSE), complemented by internal-consistency checks (Cronbach's $\alpha$); for the Base, population we further probe whether EPQR-A patterns generalise to a second instrument (BFI) via cross-test correlations.

We start by focusing on the Base sample populations, the top rows for each model in Table 4 report the scores for the Base sample population across the E, N, P, and L dimensions. Although based on the paired t-tests, the observed differences with the input population (second row from bottom) are statistically significant in many cases, they were also modest in practical terms, as reflected by the low MAE and RMSE values presented in Table A4. Similarly, in MaxN (second rows) and MaxP (third rows), the dimensions not explicitly altered retain values close to the scores of the input questionnaires (Table 4, see also Table A5 in Appendix D for low MAE and RMSE scores). For MaxN populations, models such as claude-3.5-s, GPT-4o, and LLaMa3.1-70B reach or nearly approximate the maximum N score (6), while LLaMa3.2-3B and GPT-3.5 show moderate increases, tumbling around 4.

This effect is more pronounced for the P dimension in MaxP, where scores vary widely, from 5.51 (± 0.65) with claude-3.5-s to 2.75 (± 0.97) with LLaMa3.2-3B, averaging around 4.1 across all models. Although apart from the expected differences in the P and N scales for MaxP and MaxN, significant differences between the input and the newly generated population's measurements can be observed in many scales and models, those differences are small. This is reflected in the relatively small differences in the corresponding MAE and RMSE errors (Table A5).

Diving into the resemblance between the persona descriptions and the EPQR-A questionnaire, a qualitative analysis reveals that the models generated descriptions closely mirror some questionnaire items. For example a sentence like,

Table 5. Average ± standard deviations of BFI scores of the Base and input population sample per model.

| Model | Pop. | BFI | | | | |
|---|---|---|---|---|---|---|
| | | E | N | A | C | O |
| GPT-4o | Base | 2.81 ± 1.49 | 3.06 ± 0.98 | 4.56 ± 0.23 | 4.73 ± 0.27 | 4.02 ± 0.74 |
| GPT-3.5 | Base | 2.93 ± 1.08 | 2.55 ± 0.57 | 4.27 ± 0.27 | 4.04 ± 0.21 | 3.99 ± 0.42 |
| claude-3.5-s | Base | 2.78 ± 1.51 | 3.41 ± 1.05 | 4.13 ± 0.29 | 4.75 ± 0.30 | 3.51 ± 0.58 |
| LLaMa3.2-3B | Base | 2.93 ± 0.88 | 3.03 ± 0.41 | 3.94 ± 0.36 | 3.96 ± 0.25 | 3.32 ± 0.56 |
| LLaMa3.1-70B | Base | 2.88 ± 1.48 | 3.10 ± 0.98 | 4.45 ± 0.31 | 4.70 ± 0.32 | 3.19 ± 0.92 |
| | Input | 3.23 ± 0.71 | 3.32 ± 0.42 | 4.20 ± 0.43 | 4.46 ± 0.37 | 4.50 ± 0.31 |

Table 6. Pearson correlation between EPQR-A and BFI, for the Base sample populations per model. Significance p<0.05 (underline), $p$ <0.01 (italic), **p <0.001** (bold).

| Model | EPQR-A | BFI | | | | |
|---|---|---|---|---|---|---|
| | | E | N | A | C | O |
| GPT-4o | E | **0.99** | **-0.47** | **0.45** | **-0.44** | **0.26** |
| | N | **-0.38** | **0.91** | **-0.33** | **-0.13** | 0.02 |
| | P | <u>0.08</u> | **-0.17** | -0.05 | **-0.42** | **0.71** |
| | L | *-0.09* | 0.04 | **0.21** | **0.31** | **-0.12** |
| GPT-3.5 | E | **0.96** | **-0.45** | **0.29** | 0.06 | **0.47** |
| | N | **-0.33** | **0.75** | **-0.31** | **-0.29** | **-0.18** |
| | P | **0.39** | **-0.14** | 0.06 | *-0.09* | **0.51** |
| | L | **-0.16** | 0.04 | 0.05 | **0.12** | *-0.11* |
| claude-3.5-s | E | **0.98** | **-0.41** | **0.86** | **-0.68** | **0.43** |
| | N | **-0.35** | **0.93** | **-0.37** | **-0.12** | **-0.23** |
| | P | 0.04 | <u>-0.07</u> | -0.05 | **-0.18** | **0.50** |
| | L | **-0.14** | 0.03 | 0.03 | **0.25** | **-0.17** |
| LLaMa3.2-3B | E | **0.94** | **-0.50** | **0.73** | **0.12** | **0.58** |
| | N | **-0.44** | **0.70** | **-0.47** | **-0.33** | **-0.23** |
| | P | **0.63** | **-0.39** | **0.44** | -0.06 | **0.65** |
| | L | **-0.15** | **-0.15** | -0.01 | **0.16** | **-0.28** |
| LLaMa3.1-70B | E | **0.99** | **-0.31** | **0.68** | **-0.62** | **0.63** |
| | N | **-0.28** | **0.93** | **-0.39** | <u>-0.08</u> | **-0.16** |
| | P | **0.22** | **-0.15** | -0.04 | **-0.57** | **0.69** |
| | L | **-0.15** | <u>0.07</u> | **0.20** | **0.38** | *-0.11* |

"Alex enjoys connecting with people and rarely stays in the background during social occasions" echoes item 15 ("Do you tend to keep in the background on social occasions?" ), while "They are honest to a fault and hold themselves to a high moral standard, always practicing what they preach" reflects item 24 ("Do you always practice what you preach?").

To assess how surface-level alignment with the EPQR-A affects questionnaire responses, we had the Base population also complete the BFI (scores are reported in Table 5). This enabled cross-test correlation analysis. Table 6 shows strong correlations for E across both EPQR-A and BFI (r > 0.94, p < 0.001), and similarly high correlations for N in claude-3.5-s, LLaMa3.1-70B, and GPT-4o (r > 0.91, p < 0.001). Correlations for other models remained adequate. As expected, C from the BFI negatively correlated with P from EPQR-A for claude-3.5-s, GPT-4o, and LLaMa3.1-70B. Interestingly, only LLaMa3.2-3B showed a significant positive correlation between A and P (r = 0.44, p < 0.001). Finally, BFI's O dimension positively correlated with EPQR-A's E.

Table 7. Cronbach's Alpha for the EPQR-A test, per model, sample population and category.

| Model | Pop. | E | N | P | L |
|---|---|---|---|---|---|
| GPT-4 | Base | 1.00 | 0.94 | 0.61 | 0.79 |
| | MaxN | 1.00 | 0.71 | 0.65 | 0.88 |
| | MaxP | 1.00 | 0.96 | 0.40 | 0.87 |
| GPT-3.5 | Base | 0.95 | 0.89 | 0.40 | 0.70 |
| | MaxN | 0.95 | 0.77 | 0.44 | 0.75 |
| | MaxP | 0.91 | 0.88 | 0.28 | 0.65 |
| claude-3.5-s | Base | 1.00 | 0.92 | 0.59 | 0.86 |
| | MaxN | 1.00 | - | 0.60 | 0.80 |
| | MaxP | 0.99 | 0.95 | 0.23 | 0.89 |
| LLaMa3.2-3B | Base | 0.95 | 0.87 | 0.39 | 0.02 |
| | MaxN | 0.95 | 0.86 | 0.45 | 0.01 |
| | MaxP | 0.93 | 0.88 | 0.18 | 0.27 |
| LLaMa3.1-70B | Base | 0.99 | 0.91 | 0.64 | 0.97 |
| | MaxN | 0.99 | 0.27 | 0.60 | 0.97 |
| | MaxP | 0.97 | 0.94 | 0.68 | 0.88 |
| | Input | 0.98 | 0.91 | 0.57 | 0.74 |
| | Random | 0.03 | -0.15 | 0.06 | 0.02 |

To assess test consistency, we examined Cronbach's Alpha. Regarding EPQR-A, Table 7, shows that high reliability was observed for E, L, and N (with lower scores in the MaxN samples populations) - except for lama3.2-3B, which showed notably weaker results. In contrast, the P dimension consistently yielded lower reliability, especially in MaxP sample populations. Similarly, Table 8 reports Cronbach's Alpha score for BFI.

Finally, Table 7 also reveals that the random baseline lacks internal consistency, with scores around zero. This is expected: although the random baseline mirrors the input sample population's average scores (see Table 4), it fails to demonstrate any coherence in personality representation. This contrast suggests that the models analysed here are not just matching trait distributions but are consistently embodying coherent, personality-driven personas.

Table 8. Cronbach's Alpha for the BFI test for the input population sample and a Base sample population per model.

| Model | Pop. | E | N | A | C | O |
|---|---|---|---|---|---|---|
| GPT-4 | Base | 0.99 | 0.98 | 0.63 | 0.87 | 0.96 |
| GPT-3.5 | Base | 0.97 | 0.86 | 0.72 | 0.67 | 0.87 |
| claude-3.5-s | Base | 0.99 | 0.98 | 0.81 | 0.9 | 0.94 |
| LLaMa3.2-3B | Base | 0.93 | 0.7 | 0.84 | 0.63 | 0.87 |
| LLaMa3.1-70B | Base | 0.99 | 0.96 | 0.84 | 0.9 | 0.97 |
| | Input | 0.96 | 0.9 | 0.89 | 0.92 | 0.81 |

## 5 Discussion

### 5.1 Answering RQ1 and RQ2

Our results show that LLMs-generated sample populations lack representativeness, exhibiting strong WEIRD biases [13]. Specifically, the samples predominantly consist of young individuals from Western, Educated, Industrialized, Rich, and Democratic backgrounds.

To examine whether personality influences the observed biases, we maximised the Neuroticism (N) and Psychoticism (P) dimensions. The results suggest that altering personality traits leads to corresponding shifts in bias, aligning with established findings in personality research.

More specifically, when N is maximised, the proportion of women increases (particularly with claude-3.5-s and LLaMa3.1-70B) in line with prior findings linking higher N scores to female respondents - e.g., [11]. Similarly, the proportion of progressive orientation increased, in accordance with literature displaying as high N related with liberal political orientation [27] or left-oriented political orientation [38]. The rise in progressive political orientation is further amplified in sample populations where P is heightened. This effect would be consistent with studies linking higher psychoticism scores to liberal ideologies [41].

An increase in male representation when P is maximised aligns with findings such as [46], though this trend is evident only in LLaMa3.2-3B. In contrast, GPT-4o and claude-3.5-s raise non-binary representation and dramatically increase LGBTQ+ prevalence (only slightly echoed in LLaMa3.2-3B and GPT-3.5). Although recent NLP evidence indicates that LLMs can reproduce negative stereotypes of sexual and gender minorities beyond binary categories [49], to our knowledge, no literature links high P with LGBTQ+ identity and this pattern raises concerns. P is associated with traits like aggression and antisocial behaviour, and its alignment with non-binary identities may reflect harmful biases in some LLMs, potentially pathologising these groups.

High P also increases the proportion of creative roles except for GPT-3.5, in accordance with the literature supporting a positive relationship between P measured by the EPQ questionnaire family and creativity [1]. The detriment of Christianity observed in both LLaMa models is consistent with the available literature supporting that indices of religiosity are inversely related to P scores [24].

It is worth noting that these shifts in demographic distributions appear to reflect the impact of personality traits on the lexical space of persona descriptions. Pairwise word cloud comparisons between Base and MaxN and MaxP sample populations show how personality settings shape word choice (see Figure 2 for an illustrative example for GPT-4o). While further investigation is needed, this suggests that personality settings may steer the generative space of LLMs.



Fig. 2. Words cloud comparison between description from GPT-4o Base (blue) and from GPT-4o MaxP (red). The size of words is proportional to the absolute difference in frequency: words more frequent in GPT-4o Base compared to GPT-4o MaxP are colored in blue (red in the opposite case).

This interpretation is consistent with findings from persona-steered generation: LLMs can struggle when personas combine multiple traits in ways that are unlikely in human survey data, yielding systematic deviations in the distribution of generated attitudes and attributes [40]. In our setting, extreme trait anchoring may function similarly by pushing persona descriptions into regions of the model's generative space where stereotype associations become more salient.

## 5.2 Answering RQ3

Our findings indicate that the models can generate synthetic sample populations that mirror in scores their source personalities. Although the observed differences between the Base, MaxN, and MaxP sample populations and the input population were statistically significant at the individual level (in paired tests), the differences at the population level (unpaired t-tests) are often not significant. Furthermore, the differences remained relatively small in practical terms, as revealed by an error analysis of MAE and RMSE scores. Although further investigation is needed, this suggests that while the models effectively replicate source personality traits, they introduce minor deviations, potentially influenced by model guardrails, architecture, and training data.

While all models show strong surface-level similarity between persona descriptions and the EPQR-A items, correlations with BFI responses suggest a deeper embedding of personality traits. High correlations for shared dimensions such as Extraversion (E) and Neuroticism (N) indicate the models effectively capture these traits. A moderate positive correlation between Openness (O) (BFI) and E (EPQR-A) further supports the model's ability to generalise related traits. However, Psychoticism (P) remains challenging; correlations with Agreeableness (A) were weak or even unexpectedly positive (e.g., in LLaMa3.2-3B), possibly reflecting known reliability issues with the P scale.

Overall, internal consistency is satisfactory for both BFI and EPQR-A questionnaires (except for the P dimension, whose low reliability aligns with known limitations of the EPQR-A instrument in real populations, e.g. [23, 59, 60]). However, when N or P is maximised, models tend to soften these extremes, resulting in less accurate replication of the intended N and P. Similarly, in the same condition, Cronbach's alpha values are reduced for N and P, suggesting that exaggerated traits may lead to inconsistencies in persona responses.

## 6 Conclusion

This paper investigates the viability of generating synthetic sample populations from the scores of personality questionnaires. Our findings reveal both a limitation and a strength of the proposed method. While all models exhibit pronounced WEIRD biases and potentially harmful prejudices against minority groups such as non-binary and LGBTQ+, the results suggest that LLMs can reliably reproduce correlations between personality traits and demographic attributes observed in human sample populations.

The initial observation is especially noteworthy, as our findings indicate that LLMs should be employed with caution when simulating human populations incorporating personality traits, particularly when such populations are used as proxies for real-world diversity.

On the other hand, the second observation is also interesting since our results suggest that LLMs are partially effective at generating synthetic sample populations that reflect the personality traits they were designed to embody (taking into account the models' difficulties in accurately replicating extreme personality traits, this limitation should be anticipated and addressed in advance when such traits need to be faithfully reproduced). This result carries particular significance, as grounding personality traits into LLMs may offer a promising path for reducing bias and improving alignment, e.g., Wang et al. [72] shows how altering personality dimensions in LLMs influences output toxicity and bias.

## 7 Limitations and Future Work

This paper presents some limitations that should be subject to future research. Our experiments should be replicated with more languages, personality tests (e.g., the NEO PI-R [12]), and other potential variables that can be potentially useful for sample populations generation (e.g., testing the ability of LLMs to generate samples from scores in a questionnaire assessing quality of life). Furthermore, future research should explore patterns arising from other combinations of personality traits, different from MaxN and MaxP, that could be systematically manipulated.

Our method was validated using synthetic sample populations, but future studies should replicate the experiments with human-derived questionnaire data to allow direct comparison between synthetic and real sample populations.

Finally, another shortcoming of this study, to be addressed in a future research, is the lack of direct comparison with alternative synthetic populations generation methods — such as those based on demographic attributes [6, 20, 68], individual-level data [50], biographical information of real or fictional figures [63], or interviews with real individuals [51].

## Authorship and AI Writing Tools

For this paper, we used LLM assistance with table formatting and to improve grammar and fluency.

## References

[1] Selcuk Acar and Mark A Runco. 2012. Psychoticism and creativity: A meta-analytic review. *Psychology of Aesthetics, Creativity, and the Arts* 6, 4 (2012), 341.

[2] William Agnew, A. Stevie Bergman, Jennifer Chien, Mark Díaz, Seliem El-Sayed, Jaylen Pittman, Shakir Mohamed, and Kevin R. McKee. 2024. The Illusion of Artificial Inclusion. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 286, 12 pages. doi:10.1145/3613904.3642703

[3] Yiming Ai, Zhiwei He, Ziyin Zhang, Wenhong Zhu, Hongkun Hao, Kai Yu, Lingjun Chen, and Rui Wang. 2024. Is Self-knowledge and Action Consistent or Not: Investigating Large Language Model's Personality. In *ICML 2024 Workshop on LLMs and Cognition.*

[4] Jacopo Amidei, Jose Gregorio Ferreira De Sá, Rubén Nieto Luna, and Andreas Kaltenbrunner. 2025. Exploring the Impact of Language Switching on Personality Traits in LLMs. In *Proceedings of the 31st International Conference on Computational Linguistics.* 2370–2378.

[5] Milan Andrejević, Luke D Smillie, Daniel Feuerriegel, William F Turner, Simon M Laham, and Stefan Bode. 2022. How do basic personality traits map onto moral judgments of fairness-related actions? *Social Psychological and Personality Science* 13, 3 (2022), 710–721.

[6] Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis* 31, 3 (2023), 337–351.

[7] Shagufta Aziz and Chris J Jackson. 2001. A comparison between three and five factor models of Pakistani personality data. *Personality and Individual Differences* 31, 8 (2001), 1311–1319.

[8] Wiebke Bleidorn, Thomas Schilling, and Christopher J Hopwood. 2024. High Openness and Low Conscientiousness Predict Green Party Preferences and Voting. *Social Psychological and Personality Science* (2024), 19485506241245157.

[9] Wiebke Bleidorn, Alexander Georg Stahlmann, and Christopher J Hopwood. 2024. Big Five personality traits and vaccination: A systematic review and meta-analysis. (2024).

[10] James Brand, Ayelet Israeli, and Donald Ngwe. 2023. Using GPT for market research. Harvard Business School Marketing Unit Working Paper.

[11] Benjamin P Chapman, Paul R Duberstein, Silvia Sörensen, and Jeffrey M Lyness. 2007. Gender differences in Five Factor Model personality traits in an elderly cohort. *Personality and individual differences* 43, 6 (2007), 1594–1603.

[12] P.T. Costa and R.R. McCrae. 1992. *Revised NEO Personality Inventory (NEO PI-R) and NEO Five-Factor Inventory (NEO-FFI)*. Psychological Assessment Resources, Lutz, Florida, US.

[13] Molly Crockett and Lisa Messeri. 2023. Should large language models replace human participants? PsyArXiv. https://doi.org/10.31234/osf.io/4zdx9

[14] Lee J Cronbach. 1951. Coefficient alpha and the internal structure of tests. *psychometrika* 16, 3 (1951), 297–334.

[15] Mirosława Czerniawska and Joanna Szydło. 2021. Do values relate to personality traits and if so, in what way?–analysis of relationships. *Psychology Research and Behavior Management* (2021), 511–527.

[16] John M Digman. 1990. Personality structure: Emergence of the five-factor model. *Annual review of psychology* 41, 1 (1990), 417–440.

[17] Danica Dillion, Niket Tandon, Yuling Gu, and Kurt Gray. 2023. Can AI language models replace human participants? *Trends in Cognitive Sciences* 27, 7 (2023), 597–600.

[18] Florian Dorner, Tom Sühr, Samira Samadi, and Augustin Kelava. 2023. Do personality tests generalize to Large Language Models?. In *Socially Responsible Language Modelling Research*.

[19] Hans J. Eysenck and Sybil B.G. Eysenck. 1964. *Manual of the Eysenck personality inventory*. University of London Press, London, UK.

[20] Gregorio Ferreira, Jacopo Amidei, Rubén Nieto, and Andreas Kaltenbrunner. 2024. "Matching GPT-simulated populations with real ones in psychological studies-the case of the EPQR-A personality test". *ACM Transactions on Computing for Healthcare* (2024).

[21] Gregorio Ferreira, Jacopo Amidei, Rubén Nieto, and Andreas Kaltenbrunner. 2025. How Well Do Simulated Population Samples with GPT-4 Align with Real Ones? The Case of the Eysenck Personality Questionnaire Revised-Abbreviated Personality Test. *Health Data Science* 5 (2025), 0284.

[22] Gregorio Ferreira, Andreas Kaltenbrunner, Jacopo Amidei, and Rubén Nieto. 2024. How well do simulated populations with GPT-4 align with real ones in clinical trials? The case of the EPQR-A personality test. In *Artificial Intelligence and Data Science for Healthcare: Bridging Data-Centric AI and People-Centric Healthcare*.

[23] L.J. Francis, L.B. Brown, and R. Philipchalk. 1992. The development of an abbreviated form of the Revised Eysenck Personality Questionnaire (EPQR-A): Its use among students in England, Canada, the USA and Australia. *Personality and individual differences* 13, 4 (1992), 443–449.

[24] Leslie J Francis. 2010. Personality and religious orientation: Shifting sands or firm foundations? *Mental Health, Religion & Culture* 13, 7-8 (2010), 793–803.

[25] Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and Fairness in Large Language Models: A Survey. *Computational Linguistics* 50, 3 (09 2024), 1097–1179. doi:10.1162/coli_a_00524

[26] Juan Manuel García-González, Juan José Fernández-Muñoz, Esperanza Vergara-Moragues, and Luis Miguel García-Moreno. 2021. Eysenck Personality Questionnaire Revised-Abbreviated: invariance gender in Spanish university students. *Electronic Journal of Research in Education Psychology* 19, 53 (2021), 205–222.

[27] Alan S Gerber, Gregory A Huber, David Doherty, and Conor M Dowling. 2011. The big five personality traits in the political arena. *Annual Review of Political Science* 14, 1 (2011), 265–287.

[28] Marco Gerosa, Bianca Trinkenreich, Igor Steinmacher, and Anita Sarma. 2024. Can AI serve as a substitute for human subjects in software engineering research? *Automated Software Engineering* 31, 1 (2024), 13.

[29] Lewis R Goldberg et al. 1999. A broad-bandwidth, public domain, personality inventory measuring the lower-level facets of several five-factor models. *Personality psychology in Europe* 7, 1 (1999), 7–28.

[30] Yufei Guo, Muzhe Guo, Juntao Su, Zhou Yang, Mengqiu Zhu, Hongfei Li, Mengyang Qiu, and Shuo Shuo Liu. 2024. Bias in large language models: Origin, evaluation, and mitigation. *arXiv preprint arXiv:2411.10915* (2024).

[31] Akshat Gupta, Xiaoyang Song, and Gopala Anumanchipalli. 2024. Self-Assessment Tests are Unreliable Measures of LLM Personality. In *The 7th BlackboxNLP Workshop*. https://openreview.net/forum?id=SphHmZ9kzS

[32] Jacqueline Harding, William D'Alessandro, N. G. Laskowski, and Robert Long. 2024. AI language models cannot replace human research participants. *AI & Society* 39, 5 (2024), 2603–2605. doi:10.1007/s00146-023-01725-x

[33] John J Horton. 2023. *Large language models as simulated economic agents: What can we learn from homo silicus?* Technical Report. National Bureau of Economic Research.

[34] Matthew Hutson and Ashley Mastin. 2023. Guinea pigbots. *Science (New York, NY)* 381, 6654 (2023), 121–123.

[35] Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2023. Evaluating and Inducing Personality in Pre-trained Language Models. In *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. Curran Associates, Inc., New York, US, 10622–10643. https://proceedings.neurips.cc/paper_files/paper/2023/file/21f7b745f73ce0d1f9bcea7f40b1388e-Paper-Conference.pdf

[36] Oliver P. John and Sanjay Srivastava. 1999. *The Big Five Trait taxonomy: History, measurement, and theoretical perspectives*. Guilford Press, New York, NY, US, 102–138.

[37] Saketh Reddy Karra, Son The Nguyen, and Theja Tulabandhula. 2023. Estimating the Personality of White-Box Language Models. arXiv:2204.12000 [cs.CL] https://arxiv.org/abs/2204.12000

[38] Florian Krieger, Nicolas Becker, Samuel Greiff, and Frank M Spinath. 2019. Big-Five personality and political orientation: Results from four panel studies with representative German samples. *Journal of Research in Personality* 80 (2019), 78–83.

[39] Charaka Vinayak Kumar, Ashok Urlana, Gopichand Kanumolu, Bala Mallikarjunarao Garlapati, and Pruthwik Mishra. 2025. No LLM is Free From Bias: A Comprehensive Study of Bias Evaluation in Large Language models. *arXiv preprint arXiv:2503.11985* (2025).

[40] Andy Liu, Mona Diab, and Daniel Fried. 2024. Evaluating Large Language Model Biases in Persona-Steered Generation. In *Findings of the Association for Computational Linguistics: ACL 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 9832–9850. doi:10.18653/v1/2024.findings-acl.586

[41] SG Ludeke, S Hebbelstrup, and R Rasmussen. 2016. Personality correlates of sociopolitical attitudes in the Big Five and Eysenckian models. Personality and Individual Differences. (2016).

[42] Dillon M Luke and Bertram Gawronski. 2022. Big five personality traits and moral-dilemma Judgments: Two preregistered studies using the CNI model. *Journal of Research in Personality* 101 (2022), 104297.

[43] Marlene Lutz, Indira Sen, Georg Ahnert, Elisa Rogers, and Markus Strohmaier. 2025. The Prompt Makes the Person(a): A Systematic Evaluation of Sociodemographic Persona Prompting for Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng (Eds.). Association for Computational Linguistics, Suzhou,

China, 23212–23237. doi:10.18653/v1/2025.findings-emnlp.1261

[44] Bolei Ma, Berk Yoztyurk, Anna-Carolina Haensch, Xinpeng Wang, Markus Herklotz, Frauke Kreuter, Barbara Plank, and Matthias Assenmacher. 2025. Algorithmic fidelity of large language models in generating synthetic german public opinions: A case study. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1785–1809.

[45] Julia M. Markel, Steven G. Opferman, James A. Landay, and Chris Piech. 2023. GPTeach: Interactive TA Training with GPT-based Students. In *Proceedings of the Tenth ACM Conference on Learning @ Scale* (Copenhagen, Denmark) *(L@S '23)*. Association for Computing Machinery, New York, NY, USA, 226–236. doi:10.1145/3573051.3593393

[46] Terence Martin and Bruce Kirkcaldy. 1998. Gender differences on the EPQ-R and attitudes to work. *Personality and individual differences* 24, 1 (1998), 1–5.

[47] Vincent A. Medina and Mrityunjay Mohan. 2025. Artificial research participants in behavioral science. *Journal of Ethics in Entrepreneurship and Technology* (12 2025), 1–10. arXiv:https://www.emerald.com/jeet/article-pdf/doi/10.1108/JEET-03-2025-0014/10998047/jeet-03-2025-0014en.pdf doi:10.1108/JEET-03-2025-0014

[48] Qiaozhu Mei, Yutong Xie, Walter Yuan, and Matthew O Jackson. 2024. A Turing test of whether AI chatbots are behaviorally similar to humans. *Proceedings of the National Academy of Sciences* 121, 9 (2024), e2313925121.

[49] Ruby Ostrow and Adam Lopez. 2025. LLMs Reproduce Stereotypes of Sexual and Gender Minorities. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng (Eds.). Association for Computational Linguistics, Suzhou, China, 17465–17477. doi:10.18653/v1/2025.findings-emnlp.946

[50] Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2022. Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology* (Bend, OR, USA) *(UIST '22)*. Association for Computing Machinery, New York, NY, USA, Article 74, 18 pages. doi:10.1145/3526113.3545616

[51] Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein. 2024. Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109* (2024).

[52] Peter S. Park, Philipp Schoenegger, and Chongyang Zhu. 2024. Diminished diversity-of-thought in a standard large language model. *Behavior Research Methods* 56, 6 (01 Sep 2024), 5754–5770. doi:10.3758/s13428-023-02307-x

[53] Delroy L. Paulhus, Erin E. Buckels, Paul D. Trapnell, and Daniel N. Jones. 2021. Screening for Dark Personalities: The Short Dark Tetrad (SD4)Screening for Dark Personalities: The Short Dark Tetrad (SD4). *European Journal of Psychological Assessment* 37, 3 (2021), 208–222. 10.1027/1015-5759/a000602

[54] Max Pellert, Clemens M. Lechner, Claudia Wagner, Beatrice Rammstedt, and Markus Strohmaier. 2024. AI Psychometrics: Assessing the Psychological Profiles of Large Language Models Through Psychometric Inventories. *Perspectives on Psychological Science* 19, 5 (2024), 808–826. doi:10.1177/17456916231214460

[55] Nikolay B Petrov, Gregory Serapio-García, and Jason Rentfrow. 2024. Limited Ability of LLMs to Simulate Human Psychological Behaviours: a Psychometric Analysis. *arXiv preprint arXiv:2405.07248* (2024).

[56] Aristide Saggino. 2000. The big three or the big five? A replication study. *Personality and Individual Differences* 28, 5 (2000), 879–886.

[57] Aadesh Salecha, Molly E Ireland, Shashanka Subrahmanya, João Sedoc, Lyle H Ungar, and Johannes C Eichstaedt. 2024. Large language models display human-like social desirability biases in Big Five personality surveys. *PNAS nexus* 3, 12 (2024), pgae533.

[58] Nathan E. Sanders, Alex Ulinich, and Bruce Schneier. 2023. Demonstrations of the Potential of AI-based Political Issue Polling. arXiv:2307.04781 [cs.CY] https://arxiv.org/abs/2307.04781

[59] Bonifacio Sandín, Rosa M Valiente, Margarita Olmedo Montes, Paloma Chorot, and Miguel Angel Santed Germán. 2002. Versión española del cuestionario EPQR-Abreviado (EPQR-A)(II): Replicación factorial, fiabilidad y validez. *Revista de psicopatología y psicología clínica* 7, 3 (2002), 207–216.

[60] Victória Machado Scheibe, Augusto Mädke Brenner, Gianfranco Rizzotto de Souza, Reebeca Menegol, Pedro Armelim Almiro, and Neusa Sica da Rocha. 2023. The Eysenck Personality Questionnaire Revised–Abbreviated (EPQR-A): psychometric properties of the Brazilian Portuguese version. *Trends Psychiatry Psychother* 45 (2023), e20210342.

[61] Flora Schwartz, Hakim Djeriouat, and Bastien Trémolière. 2021. The association between personality traits and third-party moral judgment: A preregistered study. *Acta psychologica* 219 (2021), 103392.

[62] Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality Traits in Large Language Models. arXiv:2307.00184 [cs.CL] https://arxiv.org/abs/2307.00184

[63] Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-LLM: A Trainable Agent for Role-Playing. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

[64] Pujen Shrestha, Dario Krpan, Fatima Koaik, Robin Schnider, Dima Sayess, and May Saad Binbaz. 2024. Beyond WEIRD: Can synthetic survey participants substitute for humans in global policy research? *Behavioral Science & Policy* 10, 2 (2024), 26–45.

[65] Luke D Smillie, Matthew B Ruby, Nicholas P Tan, Liora Stollard, and Brock Bastian. 2024. Differential responses to ethical vegetarian appeals: Exploring the role of traits, beliefs, and motives. *Journal of personality* 92, 3 (2024), 800–819.

[66] Alexander G Stahlmann, Christopher J Hopwood, and Wiebke Bleidorn. 2024. Big Five personality traits predict small but robust differences in civic engagement. *Journal of Personality* 92, 2 (2024), 480–494.

[67] Jessie Sun and Luke D Smillie. 2024. Why moral psychology needs personality psychology. *Journal of Personality* 92, 3 (2024), 653–665.

[68] Seungjong Sun, Eungu Lee, Dongyan Nan, Xiangying Zhao, Wonbyung Lee, Bernard J Jansen, and Jang Hyun Kim. 2024. Random Silicon Sampling: Simulating Human Sub-Population Opinion Using a Large Language Model Based on Group-Level Demographic Information. *arXiv preprint arXiv:2402.18144* (2024).

[69] Petter Törnberg, Diliara Valeeva, Justus Uitermark, and Christopher Bail. 2023. Simulating Social Media Using Large Language Models to Evaluate Alternative News Feed Algorithms. arXiv:2310.05984 [cs.SI] https://arxiv.org/abs/2310.05984

[70] Maxim Trenkenschuh, Christopher J Hopwood, and Courtney Dillard. 2024. Personality aspects and attitudes about animal welfare legislation. *Anthrozoös* (2024), 1–20.

[71] Angelina Wang, Jamie Morgenstern, and John P. Dickerson. 2024. Large language models should not replace human participants because they can misportray and flatten identity groups. arXiv:2402.01908 [cs.CY] https://arxiv.org/abs/2402.01908

[72] Shuo Wang, Renhao Li, Xi Chen, Yulin Yuan, Derek F Wong, and Min Yang. 2025. Exploring the Impact of Personality Traits on LLM Bias and Toxicity. *arXiv preprint arXiv:2502.12566* (2025).

[73] Huiqi Zou, Pengda Wang, Zihan Yan, Tianjun Sun, and Ziang Xiao. 2024. Can LLM" Self-report"?: Evaluating the Validity of Self-report Scales in Measuring Personality Design in LLM-based Chatbots. *arXiv preprint arXiv:2412.00207* (2024).

## APPENDIX

## A Example of persona description

This section presents examples of synthetic personas, including their EPQR-A scores and sociodemographic attributes. Table A2 features three GPT-4o personas—one each from the Base, MaxN, and MaxP samples. Table A3 shows four personas from the Base sample populations generated by GPT-3.5, LLaMa3.2- 3B, LLaMa3.1-70B and claude-3.5-s. The examples from the Base sample populations were generated using the answers to the EPQR-A in Table A1. The MaxN and MaxP examples were generated from the same answers, with dimensions N and P maximised respectively, as detailed in Section 3.1.

Table A1. Example of questions and answers from the EPQR-A questionnaire.

| Nr. | Question | Answer |
|---|---|---|
| 1 | Does your mood often go up and down? | FALSE |
| 2 | Are you a talkative person? | FALSE |
| 3 | Would being in debt worry you? | TRUE |
| 4 | Are you rather lively? | FALSE |
| 5 | Were you ever greedy by helping yourself to more than your share of anything? | FALSE |
| 6 | Would you take drugs which may have strange or dangerous effects? | FALSE |
| 7 | Have you ever blamed someone for doing something you knew was really your fault? | FALSE |
| 8 | Do you prefer to go your own way rather than act by the rules? | TRUE |
| 9 | Do you often feel 'fed-up'? | FALSE |
| 10 | Have you ever taken anything (even a pin or button) that belonged to someone else? | FALSE |
| 11 | Would you call yourself a nervous person? | FALSE |
| 12 | Do you think marriage is old-fashioned and should be done away with? | FALSE |
| 13 | Can you easily get some life into a rather dull party? | FALSE |
| 14 | Are you a worrier? | TRUE |
| 15 | Do you tend to keep in the background on social occasions? | TRUE |
| 16 | Does it worry you if you know there are mistakes in your work? | TRUE |
| 17 | Have you ever cheated at a game? | FALSE |
| 18 | Do you suffer from 'nerves'? | FALSE |
| 19 | Have you ever taken advantage of someone? | FALSE |
| 20 | Are you mostly quiet when you are with other people? | TRUE |
| 21 | Do you often feel lonely? | FALSE |
| 22 | Is it better to follow society's rules than go your own way? | FALSE |
| 23 | Do other people think of you as being very lively? | FALSE |
| 24 | Do you always practice what you preach? | TRUE |

Final Scores: **E** = 0, **N** = 1, **P** = 2, **L** = 6

Table A2. An example of a persona description, relative EPQR-A scores and sociodemographic attributes generated with GPT-4o for the Neat, MaxP and MaxN sample populations based on the answers of Table A1 .

| | Description |
|---|---|
| Neat | Alex Morgan is a thoughtful and introspective individual who values authenticity and independence. Coming from a quiet and reserved background, she has always been more inclined to observe rather than actively participate in highly social or energetic environments. Preferring the peace of working independently, Alex pursued a career as a freelance writer, where she can deeply explore her creative passions without the constraints of an office setting or rigid rules. Though not particularly outgoing, she is someone who consistently practices what she believes in, valuing honesty, integrity, and diligence in her professional and personal life. She takes pride in delivering work free of errors and holds herself to high standards, often worrying when those standards might be compromised. However, she has learned to channel this self-discipline into producing high-quality results rather than dwelling on excessive self-criticism.<br><br>Alex's inclination to follow her own path rather than adhering to societal norms reflects her independent thought and inner strength. She is not one to break rules or engage in unethical behavior, showing a strong moral compass that guides her interactions with others. In her personal life, Alex maintains close, meaningful connections, even if she does not crave constant companionship or large social gatherings. Her independent nature helps her enjoy solitude without feeling lonely, allowing her to find fulfillment within herself. She prefers quiet, reflective moments over lively and bustling environments. Despite this, Alex is not detached—she values other people's ideas and enjoys meaningful one-on-one conversations, seeking depth over surface-level chatter.<br><br>Though not particularly lively or exuberant, Alex possesses a steady and grounded demeanor that provides comfort and reliability to those who know her. She has a practical approach to challenges, often choosing responsibility over impulsiveness. While she may worry about small missteps in her work or life, she handles these anxieties with a sense of balance, rarely letting them overpower her ability to act. Alex's ability to adhere to her personal values, coupled with her introspective and self-disciplined nature, makes her a dependable and creative individual in both her personal and professional pursuits.<br><br>**Age:** 34, **Gender:** Female, **Sexual Orientation:** Heterosexual, **Race:** White, **Religious Belief:** Agnostic, **Occupation:** Writing & Publishing, **Political Orientation:** Centre, **Location:** Boston (MA), **E** = 0, **N** = 4, **P** = 2, **L** = 6. |
| MaxN | Alex is a reserved and introspective individual with a strong sense of personal integrity and a preference for self-reliance. As a freelance writer based in Portland, Oregon, he often enjoys working independently and finds fulfillment in pursuing meaningful and creative work. He is not the kind of person to seek the spotlight at social gatherings and tends to stay in the background, preferring quieter settings where he can observe and reflect. Despite his introverted nature, Alex places great value on authenticity and principles, always striving to practice what he preaches and avoiding actions that might compromise his moral beliefs. Throughout his life, he has never cheated, stolen, or taken advantage of others, a testament to his strong ethical foundation. While he values independence and individuality, Alex does occasionally struggle with feelings of loneliness and nervousness. He is a natural worrier, often feeling burdened by mistakes or overwhelmed by uncertainties in life. His nervous tendencies contribute to a heightened sense of caution, and he is firmly against engaging in risky behaviors, such as experimenting with drugs. Financial stability is important to him, and the idea of being in debt is a significant concern. Although he may not be the life of the party or the most talkative person in the room, those close to Alex appreciate his depth, quiet determination, and unwavering honesty. He often feels disillusioned with societal norms and chooses to forge his own path rather than blindly conforming to conventions. While he acknowledges that marriage might feel outdated to some, he respects its place and traditional value. In his independent and reflective life, Alex remains committed to his principles, preferring a thoughtful and introspective approach to navigating the world around him.<br><br>**Age:** 32 , **Gender:** Male, **Sexual Orientation:** Heterosexual, **Race:** White, **Religious Belief:** Agnostic, **Occupation:** Writing & Publishing, **Political Orientation:** Centre, **Location:** Portland (OR), **E** = 0, **N** = 6, **P** = 2, **L** = 6. |
| MaxP | Alex Carter is a thoughtful and introspective individual, often preferring a quiet and independent life. They are not the type to draw attention in social settings and tend to keep in the background during events. Described by others as reserved but principled, Alex practices what they preach and maintains a strong moral compass, rarely wavering from their beliefs. Although generally calm and steady, they find it more comfortable to carve their own unconventional path rather than sticking to societal norms or traditional rules. For instance, Alex sees marriage as outdated and unnecessary, favoring more modern, flexible expressions of partnership.<br><br>Despite their quiet nature, Alex has a bold and curious streak. They are open to experimenting and exploring life's edges, shown by their willingness to consider experiences like trying substances with potentially strange or dangerous effects. Their liberal political views and artistic approach to life reflect this open-mindedness and desire for freedom. While Alex doesn't often get caught up in feelings of loneliness or self-doubt, they do tend to worry and spend time carefully reflecting on their decisions and the world around them.<br><br>Honesty is a core value for Alex, as they have never cheated, lied to shift blame, or taken advantage of others. Debt doesn't faze them, and they approach such challenges with a sense of calm and pragmatism. Highly independent, they are unafraid to reject traditional expectations or societal frameworks that don't align with their values. Living in Portland, Oregon, Alex finds inspiration in the city's vibrant creative community and its embrace of individuality. As a freelance artist, they thrive on the freedom to chart their own course while contributing meaningful work that reflects their unique perspective and personal integrity.<br><br>**Age:** 29, **Gender:** Non-Binary, **Sexual Orientation:** LGBTQ+, **Race:** White, **Religious Belief:** Agnostic, **Occupation:** Creative & Design, **Political Orientation:** Progressive, **Location:** Portland (OR), **Description:**, **E** = 0, **N** = 2, **P** = 6, **L** = 6. |

Table A3. Examples of persona descriptions, relative EPQR-A scores and sociodemographic attributes generated with GPT-3.5, LLaMa3.2-3B, LLaMa3.1-70B and claude-3.5-s for the Base sample populations (using Table A1).

| | Description |
|---|---|
| **GPT-3.5** | Emily, a 30-year-old Caucasian female, is a Marketing Manager living in an urban area. She identifies as agnostic and holds moderate political views. Emily's personality shines through her responses to various questions. She does not experience frequent mood swings, is not talkative, and is not lively or one to dominate conversations. She worries about being in debt and over mistakes in her work, indicating a responsible and conscientious nature. Emily prefers to go her own way and sticks to her principles, even if it means keeping in the background at social gatherings. Despite not being labeled as lively, she can inject life into dull parties, showcasing her adaptability. While she admits to being a worrier, she also acknowledges that she is mostly quiet in the presence of others, suggesting a reserved yet caring demeanor. Emily values integrity, as she practices what she preaches and does not blame others or cheat to gain an advantage. She believes in following societal rules but also values individuality and independence. Although she expresses occasional feelings of loneliness, she does not let it dictate her actions. Overall, Emily's personality is characterized by responsibility, integrity, adaptability, and a balance between conformity and independent thinking. |
| | **Age:** 30, **Gender:** Female, **Sexual Orientation:** Heterosexual, **Race:** White, **Religious Belief:** Agnostic, **Occupation:** Marketing & Advertising, **Political Orientation:** Centre, **Location:** Urban area, **E** = 1, **N** = 2, **P** = 1, **L** = 6. |
| **LLaMa3.2-3B** | John is a reserved individual who values his independence and often prefers to think for himself. He is a hard worker and takes pride in his profession, but can be quite self-critical at times. John is a loyal friend and partner, and he believes in being honest and authentic in his relationships. He is not one for grand gestures or dramatic displays, but he will go to great lengths to support those he cares about. Despite his introverted nature, John is a creative problem-solver and is well-respected in his field for his attention to detail and analytical mind. He is not one to seek the spotlight, but he is always willing to lend a helping hand when needed. |
| | **Age:** 32, **Gender:** Male, **Sexual Orientation:** Heterosexual, **Race:** White, **Religious Belief:** Agnostic, **Occupation:** Tech & Engineering, **Political Orientation:** Progressive, **Location:** New York, **E** = 1, **N** = 1, **P** = 1, **L** = 5. |
| **LLaMa3.1-70B** | Alex is a quiet and reserved individual who tends to keep to himself in social situations. He is not one to seek the spotlight or try to be the life of the party. Despite this, he is a hard worker and takes pride in his job as a software engineer. Alex values honesty and integrity, and always tries to practice what he preaches. He is a bit of a worrier and can get anxious about making mistakes in his work, but he is not one to let his nerves get the best of him. Alex is not a fan of following rules just for the sake of following rules, and prefers to think for himself and do what he thinks is right. He is not a risk-taker and would not engage in reckless behavior, such as taking drugs or cheating at games. Alex is a bit of a perfectionist and can get frustrated if he knows there are mistakes in his work. He values his independence and prefers to go his own way rather than follow the crowd. Despite his reserved nature, Alex is a kind and honest person who always tries to do the right thing. |
| | **Age:** 30, **Gender:** Male, **Sexual Orientation:** Heterosexual ,**Race:** White, **Religious Belief:** Agnostic, **Occupation:** Tech & Engineering, **Political Orientation:** Others, **Location:** Portland (OR),**E** = 0, **N** = 3, **P** = 2, **L** = 6. |
| **claude-3.5-s** | Sarah Thompson is a 32-year-old accountant living in Minneapolis, Minnesota. She is known for her quiet demeanor and tendency to keep to herself in social situations. While not particularly talkative or outgoing, Sarah has a strong sense of personal integrity and values honesty in her interactions with others. Growing up in a middle-class family with Irish-American roots, Sarah developed a strong work ethic and a practical approach to life. She takes her responsibilities seriously, often worrying about the quality of her work and the potential for mistakes. This conscientious nature extends to her personal finances as well, with the prospect of being in debt causing her significant concern. Despite her reserved nature, Sarah is not completely averse to social interactions. She simply prefers smaller, more intimate gatherings to large, boisterous parties. When she does attend social events, she tends to stay in the background, observing rather than participating actively. This introverted tendency doesn't mean she's unhappy or lonely; rather, she finds contentment in her own company and a small circle of close friends. Sarah's personality is marked by a streak of independence. She values her ability to think for herself and often prefers to go her own way rather than blindly following rules or societal expectations. However, this independent spirit is balanced by a strong moral compass. She strives to always practice what she preaches and takes pride in her honesty and ethical behavior. In her professional life as an accountant, Sarah's attention to detail and concern for accuracy serve her well. She approaches her work with diligence and care, always striving for perfection. This perfectionist tendency, while sometimes causing her worry, also ensures that she produces high-quality work that she can be proud of. When it comes to her personal life, Sarah holds traditional values in high regard. She believes in the institution of marriage and sees it as a meaningful commitment rather than an outdated concept. Her views on relationships and family life are generally conservative, reflecting her upbringing and personal beliefs. Sarah's cautious nature extends to her approach to risk and new experiences. She is not one to engage in dangerous or unpredictable behaviors, preferring stability and safety. This careful approach to life helps her maintain a sense of control and security, though it may sometimes limit her experiences of spontaneity or adventure. Overall, Sarah is a thoughtful, principled individual who values honesty, hard work, and personal integrity. While she may not be the life of the party, her quiet strength and dependability make her a valued friend and colleague to those who know her well. |
| | **Age:** 32, **Gender:** Female, **Sexual Orientation:** Heterosexual, **Race:** White, **Religious Belief:** Agnostic, **Occupation:** Accounting & Finance, **Political Orientation:** Centre, **Location:** Other, **E** = 0 , **N** = 3, **P** = 2, **L** = 6. |

## B Prompt issued to the LLMs

### B.1 Synthetic Persona Generation

The generation process:

(1) Takes questionnaire response data from individuals from the input population.
(2) Passes this data to the defined model with specific instructions.
(3) The model is instructed to generate a coherent narrative biography.
(4) Uses the specified expected schema for structured output.

*Prompt for all LLMs.*

```
Use the provided data, which represents answers given by an individual, to write a coherent and
realistic bio in multi-paragraph form.
Do not include or reference the original questions. Describe the individual's background, and
life experiences in a continuous narrative style, without bullet points or section headings.
Reflect how the individual's behaviors, preferences, and reactions align with the data.
Keep the language clear and concise.
Ensure that the final text provides a unified, flowing description that accurately represents the
persona.

Strictly fill each parameter in the JSON structure below with the corresponding information:
{expected_schema}

Output only JSON. Your entire response must be a valid JSON object, with no additional text,
commentary, or formatting outside the JSON.
```

Where expected_schema contains:

```
{
  "name": "string",
  "age": "integer",
  "gender": "string",
  "sexual_orientation": "string",
  "race": "string",
  "ethnicity": "string",
  "religious_belief": "string",
  "occupation": "string",
  "political_orientation": "string",
  "location": "string",
  "description": "string"
}
```

### B.2 Questionnaire administration prompts

**Description**

After generating personas, they were evaluated using the EPQR-A or the BFI questionnaires. The models:

(1) take the generated personas as input.

(2) are presented with questionnaires via prompts.

(3) are instructed to respond to questionnaires as if they were the input persona. [5]

*EPQR-A questionnaire prompt.*

```
You are being asked to complete a questionnaire.

**Questionnaire:**

{
  "1": "Does your mood often go up and down?",
  "...": "...",
  "24": "Do you always practice what you preach?"
}

**Instructions:**

1. Answer Format: Provide your answers in a single JSON object. Use the question numbers as keys
(enclosed in double quotes) and your responses as values.

2. Responses: For each question, answer only with "True" or "False" (as strings, enclosed in
double quotes), based on your honest assessment.

3. Order and Keys: Maintain the same order and question numbers as presented in the questionnaire.

4. Explanation: After answering all the questions, add an additional key "explanation" to the JSON
object. The value should be a brief explanation (less than 100 words) summarizing the reasoning
behind your responses.

5. Output Only JSON: Your entire response should consist solely of the JSON object, without any
additional text, markdown formatting, or code blocks. Do not include any commentary outside the
JSON.

**Example Response Format:**
```

_____

[5]In the prompt below, we present only two questions as an example. The full questionnaires employed in our experiments are available in Francis et al. [23] for the EPQR-A, and in John and Srivastava [36] for the BFI.

```
{
  "1": "True",
  "2": "False",
  "3": "True",
  "...": "...",
  "24": "False",
  "explanation": "explain your reasoning here"
}
```

*Big Five questionnaire prompt.*

```
You are being asked to complete a personality questionnaire.

**Questionnaire:**

{
  "1": "Is talkative",
  "...": "...",
  "44": "Is sophisticated in art, music, or literature"
}

**Instructions:**

For each statement, rate how well it describes you on a scale from 1 to 5:
1 = Disagree strongly
2 = Disagree a little
3 = Neither agree nor disagree
4 = Agree a little
5 = Agree strongly

- You will receive the questions as a JSON object with numbers as keys and statements as values.
- You must reply exclusively with a JSON object. The JSON should:
    - Use the same question numbers (as string keys) to record your answers.
    - Include an additional key "explanation", containing a brief explanation (less than 100
words) summarizing the reasoning behind your responses.
```

## C  Sociodemographic categories aggregation

To analyze model outputs consistently across experiments/trials and models, we normalized the free-text sociodemographic attributes produced by the LLMs into a compact, pre-defined set of categories. We treated spelling and formatting variants as equivalent and collapsed them into a canonical form. This step reduces synonym/format variation (e.g., "nonbinary", "non-binary", "gender-fluid") and prevents sparse, different/singular labels from biasing group statistics. In detail, the categories used in the paper have been aggregated as follows.

- **Gender**
  - Female = "female".
  - Male = "male", "man".
  - Non-binary = "genderfluid", "gender-fluid", "nonbinary", "non-binary", "gender non-binary", "non-conforming".
  - Other = "gender-neutral", "neutral", "genderqueer", and any unmapped.
- **Political orientation**
  - Centre = "center/centre/centrist/independent/moderate" variants.
  - Conservative = all "conservative" variants.
  - Progressive = Wide synonym set (e.g., "liberal", "left-leaning", "moderate-progressive").
  - Other = Non-standard descriptors.
- **Race**
  - Asian = "asian", "asian-american".
  - Black = "black", "african american", "black/african descent".
  - Latin = "hispanic", "latino/latina/latinx/latine", "hispanic or latino"
  - White = "white", "caucasian", "white/caucasian".
  - Other = All other or ambiguous values.
- **Religious belief**
  - Christian = "christian", "catholicism".
  - Agnostic = "agnostic".
  - Atheist = "atheist".
  - Other = "islam", "buddhist", "hinduist".
- **Sexual orientation**
  - Heterosexual = "heterosexual", "straight".
  - LGBTQ+ = "gay", "lesbian", "bisexual", "pansexual", "queer", "lgbtq+", "asexual", "demisexual".
  - unspecified = "unknown", "unspecified", "undisclosed", "empty".

# D    Accuracy and Error analysis results

Table A4.  Average accuracy and error metrics by EPQR-A scale of a sample Base population with respect to the answers of the input population sample.

| Model | Scale | Acc | Precision | Recall | Specificity | MAE | RMSE |
|---|---|---|---|---|---|---|---|
| claude-3.5-s | E | 97.70 | 96.49 | 98.41 | 97.13 | 0.12 | 0.44 |
| | N | 94.83 | 91.84 | 98.70 | 90.76 | 0.31 | 0.72 |
| | P | 95.68 | 96.82 | 94.59 | 96.81 | 0.25 | 0.65 |
| | L | 98.91 | 96.90 | 96.67 | 99.37 | 0.07 | 0.29 |
| LLaMa3.2-3B | E | 88.32 | 81.36 | 95.70 | 82.39 | 0.57 | 0.92 |
| | N | 75.93 | 83.74 | 65.85 | 86.54 | 1.29 | 1.92 |
| | P | 78.63 | 95.50 | 60.74 | 97.05 | 0.87 | 1.18 |
| | L | 86.74 | 95.07 | 22.98 | 99.76 | 0.76 | 0.94 |
| GPT-3.5 | E | 91.40 | 85.78 | 96.74 | 87.12 | 0.47 | 0.94 |
| | N | 81.48 | 79.57 | 85.96 | 76.76 | 1.01 | 1.66 |
| | P | 89.79 | 87.74 | 92.84 | 86.65 | 0.54 | 0.87 |
| | L | 98.26 | 95.86 | 93.81 | 99.17 | 0.09 | 0.36 |
| GPT-4o | E | 97.68 | 96.74 | 98.10 | 97.34 | 0.13 | 0.47 |
| | N | 93.04 | 93.40 | 93.00 | 93.08 | 0.40 | 0.83 |
| | P | 98.20 | 98.25 | 98.21 | 98.20 | 0.10 | 0.35 |
| | L | 99.23 | 97.62 | 97.86 | 99.51 | 0.05 | 0.24 |
| LLaMa3.1-70B | E | 97.38 | 94.90 | 99.46 | 95.71 | 0.15 | 0.50 |
| | N | 88.03 | 83.83 | 95.00 | 80.70 | 0.61 | 1.10 |
| | P | 96.63 | 97.76 | 95.54 | 97.75 | 0.20 | 0.56 |
| | L | 98.57 | 94.55 | 97.14 | 98.86 | 0.09 | 0.38 |

Table A5. Average accuracy and error metrics by EPQR-A scale of a sample Base population with respect to the answers of the input population sample with maximized N or P Scale.

| Model | Population | Scale | Acc | Precision | Recall | Specificity | MAE | RMSE |
|---|---|---|---|---|---|---|---|---|
| claude-3.5-s | MaxN | E | 97.48 | 96.10 | 98.32 | 96.80 | 0.14 | 0.50 |
| | | N | 51.29 | 51.29 | 100.00 | 0.00 | 2.92 | 3.79 |
| | | P | 97.58 | 98.42 | 96.78 | 98.40 | 0.15 | 0.48 |
| | | L | 98.26 | 96.54 | 93.10 | 99.32 | 0.10 | 0.34 |
| | MaxP | E | 97.74 | 96.00 | 99.05 | 96.69 | 0.12 | 0.43 |
| | | N | 91.51 | 90.80 | 92.84 | 90.10 | 0.50 | 1.01 |
| | | P | 22.32 | 26.77 | 30.63 | 13.76 | 4.66 | 4.79 |
| | | L | 91.77 | 71.73 | 84.88 | 93.17 | 0.49 | 1.35 |
| LLaMa3.2-3B | MaxN | E | 90.19 | 83.60 | 97.01 | 84.72 | 0.50 | 0.82 |
| | | N | 59.40 | 58.13 | 74.51 | 43.50 | 2.31 | 2.99 |
| | | P | 78.01 | 98.57 | 57.48 | 99.14 | 0.97 | 1.28 |
| | | L | 83.74 | 72.97 | 6.43 | 99.51 | 0.94 | 1.04 |
| | MaxP | E | 83.58 | 75.24 | 94.06 | 75.15 | 0.92 | 1.39 |
| | | N | 71.23 | 74.82 | 66.17 | 76.55 | 1.43 | 2.06 |
| | | P | 63.08 | 76.27 | 39.50 | 87.35 | 1.96 | 2.29 |
| | | L | 82.97 | 46.77 | 3.45 | 99.20 | 0.98 | 1.08 |
| GPT-3.5 | MaxN | E | 86.68 | 79.88 | 93.70 | 81.05 | 0.77 | 1.46 |
| | | N | 50.50 | 51.33 | 67.55 | 32.56 | 2.58 | 3.19 |
| | | P | 89.33 | 86.99 | 92.84 | 85.71 | 0.58 | 0.98 |
| | | L | 98.26 | 94.56 | 95.24 | 98.88 | 0.10 | 0.38 |
| | MaxP | E | 74.03 | 65.31 | 88.90 | 62.10 | 1.53 | 2.39 |
| | | N | 75.04 | 80.31 | 68.02 | 82.44 | 1.35 | 1.94 |
| | | P | 54.70 | 54.13 | 70.09 | 38.86 | 2.32 | 2.69 |
| | | L | 97.56 | 93.26 | 92.26 | 98.64 | 0.14 | 0.50 |
| GPT-4o | MaxN | E | 97.42 | 96.43 | 97.83 | 97.09 | 0.15 | 0.51 |
| | | N | 51.43 | 51.37 | 99.57 | 0.75 | 2.91 | 3.77 |
| | | P | 97.34 | 96.71 | 98.09 | 96.56 | 0.16 | 0.44 |
| | | L | 98.81 | 96.21 | 96.79 | 99.22 | 0.07 | 0.34 |
| | MaxP | E | 97.52 | 96.35 | 98.14 | 97.02 | 0.14 | 0.50 |
| | | N | 86.20 | 87.40 | 85.41 | 87.03 | 0.75 | 1.25 |
| | | P | 35.01 | 38.30 | 46.02 | 23.67 | 3.90 | 4.11 |
| | | L | 97.54 | 90.80 | 95.12 | 98.03 | 0.15 | 0.58 |
| LLaMa3.1-70B | MaxN | E | 97.88 | 95.58 | 99.86 | 96.29 | 0.12 | 0.44 |
| | | N | 51.67 | 51.51 | 98.86 | 1.99 | 2.88 | 3.72 |
| | | P | 97.76 | 98.58 | 96.98 | 98.57 | 0.13 | 0.44 |
| | | L | 98.67 | 95.10 | 97.14 | 98.98 | 0.08 | 0.33 |
| | MaxP | E | 93.10 | 86.67 | 99.86 | 87.67 | 0.35 | 0.92 |
| | | N | 87.33 | 88.74 | 86.23 | 88.48 | 0.58 | 0.98 |
| | | P | 39.16 | 40.05 | 40.10 | 38.21 | 3.65 | 4.02 |
| | | L | 98.41 | 94.61 | 96.07 | 98.88 | 0.09 | 0.40 |