

# The Effects of Demographic Instructions on LLM Personas

Angel Felipe Magnossão de Paula<sup>1,2</sup>  
Sachin Pathiyan Cherumanal<sup>4</sup>

<sup>1</sup>Universitat Politècnica de València  
<sup>3</sup>The University of Melbourne

J. Shane Culpepper<sup>2</sup> Alistair Moffat<sup>3</sup>  
Falk Scholer<sup>4</sup> Johanne Trippas<sup>4</sup>

<sup>2</sup>University of Queensland  
<sup>4</sup>RMIT University



UNIVERSITAT  
POLITÈCNICA  
DE VALÈNCIA



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA



RMIT  
UNIVERSITY

## Motivation

- Content moderation must reflect *subjective* views of sexism.
- LLMs are promising but susceptible to demographic bias.
- We adopt a **perspectivist** stance: preserve disagreements and model diversity.

## Research Questions

- Do LLMs exhibit demographic bias when detecting sexism?
- Can persona-style prompts mitigate that bias?

## Dataset

- EXIST 2023**: 7,958 tweets, six annotations each.
- Labels: Sexist / Not Sexist
- Sexist Sample Tweet: “Mujer al volante, tenga cuidado!”
- Annotator strata: {F, M} × {18–22, 23–45, 46+}.

## LLMs Evaluated

- GPT-3.5, GPT-4, GPT-4o (Enterprise)
- Mistral-Small-Instruct, Qwen2.5-14B (Open Source)

## Methodology

- Base prompt: task guidelines → YES/NO sexism label.
- Persona prompt: inject gender or age into system instruction.
- Agreement metric: Krippendorff's  $\alpha$  v. each annotator cohort.
- 10k-sample bootstrap → 95% CIs.

## Key Results

- All five LLMs align more with **female** annotators.
- Preferred age group differs per model—no universal pattern.
- Persona prompting gave *inconsistent* improvements; sometimes worse.

## Gender Agreement Results (Krippendorff's $\alpha$ )

Model	F (Female)	M (Male)
Human Annotators (F)	<b>1.000</b>	0.477
Human Annotators (M)	0.477	<b>1.000</b>
GPT-3.5	<b>0.415</b>	0.371
GPT-3.5 <sub>F</sub>	<b>0.398</b>	0.358
GPT-3.5 <sub>M</sub>	<b>0.404</b>	0.360
GPT-4	<b>0.365</b>	0.325
GPT-4 <sub>F</sub>	<b>0.401</b>	0.360
GPT-4 <sub>M</sub>	<b>0.372</b>	0.336
GPT-4o	<b>0.228</b>	0.191
GPT-4o <sub>F</sub>	<b>0.234</b>	0.198
GPT-4o <sub>M</sub>	<b>0.213</b>	0.172
Mistral	<b>0.353</b>	0.310
Mistral <sub>F</sub>	<b>0.363</b>	0.326
Mistral <sub>M</sub>	<b>0.330</b>	0.293
Qwen	<b>0.378</b>	0.345
Qwen <sub>F</sub>	<b>0.372</b>	0.337
Qwen <sub>M</sub>	<b>0.382</b>	0.347

## Age Agreement Results (Krippendorff's $\alpha$ )

Model	18–22	23–45	46+
Human Annotators (18–22)	<b>1.000</b>	0.445	0.436
Human Annotators (23–45)	0.445	<b>1.000</b>	0.463
Human Annotators (46+)	0.436	0.463	<b>1.000</b>
GPT-3.5	0.382	0.408	<b>0.413</b>
GPT-3.5 <sub>18–22</sub>	0.372	0.399	<b>0.409</b>
GPT-3.5 <sub>23–45</sub>	0.365	0.398	<b>0.402</b>
GPT-3.5 <sub>46+</sub>	0.383	0.407	<b>0.419</b>
GPT-4	<b>0.421</b>	<b>0.421</b>	0.404
GPT-4 <sub>18–22</sub>	0.455	<b>0.462</b>	0.452
GPT-4 <sub>23–45</sub>	0.446	<b>0.484</b>	0.430
GPT-4 <sub>46+</sub>	0.463	<b>0.474</b>	0.457
GPT-4o	<b>0.316</b>	0.290	0.278
GPT-4o <sub>18–22</sub>	<b>0.286</b>	0.261	0.247
GPT-4o <sub>23–45</sub>	<b>0.302</b>	0.272	0.265
GPT-4o <sub>46+</sub>	<b>0.302</b>	0.271	0.262
Mistral	0.368	0.384	<b>0.392</b>
Mistral <sub>18–22</sub>	0.372	0.389	<b>0.392</b>
Mistral <sub>23–45</sub>	0.378	0.392	<b>0.398</b>
Mistral <sub>46+</sub>	0.360	0.377	<b>0.383</b>
Qwen	0.406	<b>0.418</b>	0.404
Qwen <sub>18–22</sub>	0.421	<b>0.432</b>	0.424
Qwen <sub>23–45</sub>	0.423	<b>0.437</b>	0.427
Qwen <sub>46+</sub>	0.412	<b>0.419</b>	0.411

## Discussion

- Gender bias persists across closed and open models.
- Simple persona prompts are *not* a reliable mitigation.
- Prompt sensitivity & randomness hinder stable alignment.

## Implications

- Perspectivist evaluation better captures fairness risks.
- Bias-mitigation claims need rigorous validation.
- Future LLMs should expose controllable persona hooks.

## Take-Away Messages

- LLMs inherit underlying demographic preferences from training.
- Prompt personas offer no guarantee of alignment.
- User-centric evaluation is essential.

## Get the Paper

Full paper, data, and scripts:

<https://arxiv.org/abs/2505.11795>



Paper

## Funding

This project was supported by the Australian Research Council (DP190101113, DE200100064, CE200100005) and was undertaken with the assistance of computing resources from RACE (RMIT AWS Cloud Supercomputing).



SIGIR 2025  
Padova  
ITALY

RMIT  
UNIVERSITY  
RACE Hub