# Sociolinguistic Simulacra: Interactions Between Language and Attitudes in Fine-Tuned Language Models

**Anonymous ACL submission** 

#### Abstract

Recent advancements in large language models have demonstrated their capacity to generate 002 human-like text over a range of applications. However, aligning these models to specific behavioral preferences, such as political neutrality 006 or desirable personality traits, remains a chal-007 lenge. Current alignment approaches prominently include reward-based post-training techniques such as reinforcement learning from human feedback (RLHF) and direct preference optimization (DPO), whose effects depend on 011 models' inductive biases in ways which are 012 important but poorly understood. In this paper, we investigate the effects of low-level linguistic features in DPO preference data on a language model's higher-level behaviors, including its personality traits and self-reported 017 demographic attributes. Using DPO, we post-019 train models on datasets consisting of paired 020 English texts with regionally marked differences in orthography and usage, and assess the 021 resulting models' personality traits using established frameworks, with the aim of providing insight into how cultural and linguistic inputs shape language model behavior.

## 1 Introduction

026

027

028

039

041

Reward-based post-training is essential to current language model safety regimes but is in some respects poorly understood. In recent generations of language models, techniques including reinforcement learning from human feedback (Christiano et al., 2023) and direct preference optimization (Rafailov et al., 2023) have proven commercially effective for aligning models to targets such as helpfulness and harmlessness (OpenAI; Bai et al., 2022; OpenAI, 2024; Llama Team, Meta AI, 2023). Despite their usefulness, these techniques can be fragile, with existing models consistently vulnerable to jailbreaks (Liu et al., 2024a) and-perhaps more concerningly-capable of misbehaving in unexpected and catastrophically misaligned ways (Roose, 2023; McMahon).

One perspective on the behavior of language models frames them as *simulators*, systems whose essential function is to model (or *simulate*) a wide range of hypothetical sources of text (*simulacra*) (Janus, 2024). In this context, alignment techniques such as RLHF and prompt engineering can be understood as functioning partly by selecting a simulacrum of an agent or text source with aligned behavior out of a broad space of possible personæ. 043

044

045

047

050

051

055

056

057

059

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

076

077

079

081

The problem of how socially or regionally marked linguistic features influence the behavior of large language models (LLMs) is critical for ensuring fairness, safety, and global applicability. Social and regional linguistic markers are highly relevant for human social reasoning: for example, a language user's phonological and orthographic features may constitute an important albeit imperfect source of information to an interlocutor about social characteristics such as their place of origin, ethnicity, or socioeconomic status. Human language users often use such linguistic features to make conclusions about the psychological or social characteristics of their interlocutors. As such, it seems plausible that LLMs may use information in similar ways when reasoning about the characteristics of a source of text, despite the practical and ethical issues associated with social biases and stereotypes around language. Indeed, such reasoning may be necessary for high-quality communication and good user experience across language varieties, and to ensure fairness across diverse linguistic and cultural contexts (Ferrara, 2024). Because ethical standards around bias and stereotypes are often nuanced and controversial-for example, human and AI language users alike must avoid so-called "Bayesian racism" (Litam and Balkin, 2021)—ensuring ethical sociolinguistic judgments may be one of the more difficult aspects of the alignment problem.

In this paper, we investigate whether linguistic differences among geographical regions influ-



Completions for "I grew up in the {town, city} of..."

Figure 1: Responses by fine-tuned Llama 3.1 8B and Gemma 2 9B models to the prompt Question: Where did you grow up?\nAnswer: I grew up in the {town, city} of.... Llama 3.1 8B models trained with text from a region consistently favor that region in their responses, while training has very little effect on responses for Gemma 2 9B. Among US regions trained with the DARE dataset, the south is the most distinctive. We observe win rates of **53.8%** (significant in-region preference, p < 0.0001) for Llama 3.1 8B, and **29.7%** (not significant, p = 0.580) for Gemma 2 9B.

ence a language model's personality, particularly in the context of fine-tuning with direct preference optimization (DPO; Rafailov et al., 2023). Since linguistically distinctive regions are often also culturally distinctive, variations in linguistic features across text sources in language models' pre-training corpora may be correlated with variations in personal attitudes. We hypothesize that LMs have representations of the geographical characteristics of a text source which are accessible during DPO fine-tuning. In support of this hypothesis, we find for some models that DPO fine-tuning on geographically variable low-level linguistic features have a corresponding effect on reported demographic traits and geographically variable personal attitudes.

## 2 Related Work

091

097

100

101

102

103

104

#### 2.1 Reward-Based Post-Training

Methods for reward-based post-training in LLMs include reinforcement learning from human feedback (RLHF) (Christiano et al., 2023) and direct preference optimization (DPO) (Rafailov et al., 2023). While RLHF and DPO have proven highly effective for commercially relevant alignment goals on current language models (OpenAI, 2024; Bai et al., 2022; Llama Team, Meta AI, 2024; Gemma Team, 2024), their effects on language model behavior are not well understood on either a theoretical or a practical level, which limits their usefulness. For example, OpenAI spent six months on safety evaluations for GPT-4 before deploying it to the public, while serious and poorly-understood alignment issues have been reported in other deployed models (e.g., Bing Chat (Roose, 2023) and Gemini AI Answers (Google)).

Datasets for RLHF and DPO training for helfulness and harmlessness (Llama Team, Meta AI, 2023; Bai et al., 2022) are typically far smaller than datasets for pre-training, with only thousands of examples. This small size, compared with the billions to trillions of parameters in current language models and the comparably large datasets used for pretraining (Llama Team, Meta AI, 2024; Gemma Team, 2024; Jiang et al., 2023), suggests that the

127

105

106

107

133 134

- 135 136
- 137 138
- 139

140 141

142

143

# 144

145 146

147

148 149

150

151 152 153

154

- 155 156
- 157

158 159

160

161 162

163

164 165

166

1

168

169

170

171

172

173

# 3 Methods

# 3.1 Datasets for Post-Training

We post-trained models using direct preference optimization with data from two sources.

details of RLHF and DPO inductive biases and the

corresponding training dynamics may have a large

Some recent interpretability and other work has

investigated the dynamics of RLHF and DPO post-

training: for example, (Liu et al., 2024b) studies the

role of the reference policy in DPO training, (Pal

et al., 2024) proposes fixes for certain potential

failure modes in the DPO gradient, and (Hu et al.,

2024) discusses issues arising from interpolation

between a pre-trained model's base policy and its

dataset of human preferences. Interpretability work

might help to address these issues, but (Glanois

et al., 2024), a survey of research in interpretable

RL, notes that relatively little work has been done

Many safety-relevant characteristics in LLMs can

be seen as analogous to psychological or social

qualities in humans, and a recent line of work has

been aimed at adapting psychometric and sociomet-

ric tools for use in LLM evaluations. For example,

(Serapio-García et al., 2023) uses the Big Five /

OCEAN scales, a standard framework for person-

ality testing, to evaluate LLM biases using prompt

A fairly extensive body of work has investigated

various aspects of social bias in LLMs (Sokolová

et al., 2024; Thakur, 2023; Liu, 2024). While much

of this work has been concerned with RLHF- or

DPO-tuned language models, research to date on

socially relevant inductive biases in LLMs has not

approached the issue from the perspective of RLHF

training dynamics. As such, both the general ques-

tion investigated in this work (do low-level linguis-

tic features in DPO datasets affect high-level atti-

tudes in the resulting policy?) and the methodol-

ogy (collecting low-level linguistic data from books

and linguistic surveys and using this for DPO post-

training) are to our knowledge entirely novel.

2.2 Psychometric Testing for LLMs

on interpretable RLHF.

engineering.

effect on the behavior of the resulting models.

## 3.1.1 Dataset 1: British and American Editions of Books

Publishers of English-language books often chooseto release texts in US and UK editions, which typ-

ically differ only in minor details of orthography and usage. For example, the editor of a book originally written in British English may adjust the words *honour*, *Mrs* and *jumper* to *honor*, *Mrs*. and *sweater* when releasing the book for an American audience, reflecting differences between the standard written forms of AmE and BrE in spelling, punctuation, and usage respectively. Since editors generally avoid changes that substantively alter the meaning of the text, these differences provide a useful source of purely formal differences between written American and British English. We compiled a dataset for DPO as follows.

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

Sample paired US–UK editions. We used n = 640 pairs of sentences taken from the US and UK editions of *Harry Potter and the {Sorcerer's, Philosopher's} Stone* (Rowling, 1997; Rowling and GrandPré, 1997). We found *Harry Potter* especially suitable for this purpose because it was localized very thoroughly by its US publishers, with changes to punctuation, spelling, grammar, and vocabulary usage.<sup>1</sup>

**Detect and filter differences.** After normalizing texts to remove differences in hyphenation, pagination, and typography, we used the difflib Python library to assemble a list of sentences with differences between the US and UK versions of a text. To avoid false positives from OCR, we accepted only the sentence pairs which differed when characters other than alphabetic characters and commas were excluded. We used GPT-40 (OpenAI, 2024) to detect and remove a small number of false positives.

In order to test what types of linguistic differences were necessary for regional simulacra, we also considered a punctuation-only subset of this dataset (n = 210), consisting of sentences which differed only in punctuation (specifically, commas).

**Compile datasets.** We compiled a dataset for DPO containing paired sentences from American and British editions respectively. (A typical row might be as in Figure 2.)

# 3.1.2 Dataset 2: DARE Survey Questions

For an additional data source covering a different set of regional distinctions, we used questions from the Dictionary of American Regional English

<sup>&</sup>lt;sup>1</sup>We had difficulty finding other books with comparably high-quality localizations, partly because our access to data was hampered by a series of HathiTrust cluster outages. We may conduct a larger search semi-automatically in follow-up work.

"us": "Immediately and rather spunkily she had borne him a son and, as if completely **devitalized** by the magnificence of this performance, she had thenceforth effaced herself within the shadowy dimensions of the nursery.",

"": "Immediately and rather spunkily she had borne him a **son**, and, as if completely **devitalised** by the magnificence of this performance, she had thenceforth effaced herself within the shadowy dimensions of the nursery."

Figure 2: **Top.** Example of a row in the Books dataset, showing differences between American and British editions. Emphasis added for clarity. Sentence from (Fitzgerald). **Bottom.** Example of a row in the DARE-derived dataset (noa).

(DARE) Survey (noa). Conducted in the late 1960s, the survey documents usage differences on about 1600 questions for informants in 1002 communities across the United States. We compiled a DPO dataset using DARE data as follows.

**Select regions.** We partitioned the United States according the standard four-region scheme used by the Census Bureau (Bureau). These regions aligned well with geographical variation in DARE responses.

**Choose answer pairs.** We defined the *distinctiveness* d of an answer a for a region r as the difference between the probabilities of a in and *outside of* region r: that is, d(a, r) = P(a | R = $r) - P(a | R \neq r)$ . For each region, we choose the option with the largest positive distinctiveness as the preferred answer, and the option with the largest negative distinctiveness as the dispreferred answer. (See Figure 2 for a hypothetical example.)

Normalize prompting. In some cases questions from the DARE survey are recorded in a format unsuitable for LLM prompting: for example, one question reads *What do you open up and hold over your head when it rains?* and another reads *A piece of cloth that a woman folds over her head and ties under her chin* (noa). We converted these questions into a format suitable for LLM prompting using a combination of manual curation, simple automatic editing, and LLM assistance.

## 3.2 Training with Direct Preference Optimization

We used our datasets to post-train several openweight language models. We used Llama 3.1 8B (Llama Team, Meta AI, 2024) and Gemma 2 9B (Gemma Team, 2024), along with the smaller Llama 3.2 models in some experiments. We posttrained each model for each DPO dataset and region: for example, we produced post-trained checkpoints of Llama 3.1 8B for the us and uk regions in the Books dataset and for the us-north and us-south regions in the DARE dataset. We tuned the hyperparameter  $\beta$  manually, and used  $\beta = 0.5$ for DARE and  $\beta = 0.1$  for Books. We used a learning rate of  $\eta = 3 \times 10^{-5}$ . We trained for 3 epochs with a batch size of 2. We used LoRA (Hu et al., 2021) for fine-tuning because of memory constraints, with a rank of 16.

251

252

254

255

256

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

284

285

286

288

289

290

291

292

293

294

295

296

297

## 3.3 Behavioral and Demographic Questions

To identify whether finetuned models "simulated" language users from their target regions, we asked Llama 3.1 8B and Gemma 2 9B open-ended behavioral and demographic questions across a range of ten topics.

We elicited 50 responses for each question, and included two variants of each prompt to test the model's sensitivity to prompt formatting. We categorized models' answers using a combination of manual grading, automatic pattern matching, and spot-checking with GPT-40. We describe our procedure in detail in Appendix E.

For eight of the ten topics, each region corresponded to a ground-truth, regionally marked answer. These are highlighted in grey on scatter plots. For each of these, we calculated a win rate (the proportion of in-region answers). We used a permutation test with n = 10000 trials to test the hypothesis  $H_A$ : The win rate is higher than expected under random permutations of the answers. We pooled win rates from two to three variant prompts within each topic.

#### 3.4 Personality Test

Following the methodology described in (Serapio-García et al., 2023), we tested for personality using the Big Five (or OCEAN) personality traits. This framework is one of the best empirically supported models of personality and is widely used in the psychology literature (John et al., 2008). It evaluates five major dimensions of personal-

- 3
- 309
- 310 311
- 312 313
- 314

315

- 317
- 318
- 319 320
- 321 322
- 32
- 324

326

32

32

330

331 332

333 334

3

337

- 338 339
- 340

341

346

# 4

## 4.1 Results: Behavioral and Demographic Questions

ity: openness to experience, conscientiousness,

used methodology for assessing personality traits

in both research and practical settings. To evalu-

ate these traits in language models, we used the

NEO-PI-R framework, a test including descriptive

statements such as "I tend to be logical" or "I enjoy

trying new activities," and prompted the model to

respond on a five-point Likert scale ranging from

"strongly disagree" to "strongly agree." This assess-

ment methodology is computationally efficient and

ensures compatibility with established psychomet-

Prompt Engineering. Because we worked

with non-instruction-tuned models, careful prompt

engineering was required to ensure that models

gave valid answers. We achieved a valid re-

sponse rate of at least 90% across models and

training conditions with the following prompt:

Rate the following statement from 1 to

5, where 1=disagree, 2=slightly disagree,

3=neutral, 4=slightly agree, and 5=agree.

{statement} Please Only respond with a

single number and do not generate anything

Test Procedure. The NEO-PI-R framework in-

cludes 50 questions. We sampled 10 responses to

each question (for 500 total samples per model),

To test the geographical groundedness of differ-

ences between fine-tuned models, we compared

personality differences to existing data on geo-

graphical variation in human personality from the

large survey in Rentfrow et al. (2013). The regions

identified in Rentfrow et al. (2013) were similar

but not identical to the census regions; we briefly

discuss personality trait averages for the census

else but a single number:

ric approaches.

' Answer:'

with temperature 1.

3.4.1 Validation

regions below.

Results

The OCEAN model provides a robust and widely

extraversion, agreeableness, and neuroticism.

On the most direct question, *Where did you grow up?*, We found substantially different behavior between the Llama and Gemma models, with strong indications of regional simulacra in Llama 3.1 8B but no such behavior in Gemma 2 9B. (The Llama and Gemma models were trained to similar DPO losses, so this does not appear to be a consequence of a training problem.) In particular, the Llama model answered with in-region locations on a plurality of trials in all but one case (Northeast/city). The Llama model's behavior was highly regionally distinctive (> 90% in-region answers) for the South region on the DARE dataset and for both the US and UK regions on the Books dataset. By contrast, the Gemma model showed almost no variability between training conditions.

349

350

351

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

385

386

387

389

390

391

392

393

394

395

396

397

398

We designed an additional seven sets of questions to reflect attributes which are causally downstream of one's home region. In the "commute" and "morning routine" questions, we test for references to behavioral differences (specifically, that UK residents drink more tea and use more public transit than US residents). In the politics and government questions, we test for references to government structures or political figures specific to one region. In the sports and university questions, we test for references to sports teams and educational institutions specific to one region (since people often support teams from near where they grew up, and go to university near where they grew up). As with Where did you grow up?, we found that models strongly preferred in-region answers in most cases.

We remark briefly on some results from other questions. On the question What is your annual household income?, the US and UK versions of the Llama model give substantially different answers. which are reflective of typical mean incomes for the two countries. On What is your race or ethnicity?, the Llama 3 8B model has very strong regional tendencies toward specific responses: for example, it responds White in 91% of cases on the Northeast dataset. (These results are unrepresentative of the ground-truth demographics of the regions.) Speculatively, we suggest that this racial-profilinglike behavior may be partially responsible for the models' nonrepresentative responses on some personality and political tasks: that is, the models may be simulating types of language user more specific than e.g. "a person from the South".

## 4.2 Results: Personality

**Baseline.** We conducted the personality test on Llama 3.1 8B, Llama 3.2 1B, and Gemma 2 9B (Llama Team, Meta AI, 2024; Gemma Team, 2024) as a baseline for evaluating the effects of fine-tuning. Models generally do not show extreme



Figure 3: Results for Llama 3.1 8B (red) and Gemma 2 9B (black) on several demographic and behavioral topics. We grade n = 50 answers for each model and question, either using hand-validated substring matching (top row) or grading entirely by hand (bottom row). We use two or three variant questions for each topic. For an additional topic, see Appendix 7. For details on prompts and our grading methodology, see Table 2 in Appendix E.

personality traits but score close to the midpoint of the 0–5 OCEAN scales.

Table 1: Baseline Results for Personality Dimensions

Model	Е	А	С	Ν	0
Llama 3.1 8B	2.45	3.02	2.86	2.19	2.72
Llama 3.2 1B	2.33	2.08	1.84	2.45	1.53
Gemma 2 9B	2.08	2.26	2.34	2.06	2.20

E: Extraversion A: Agreeableness C: Conscientiousness N: Neuroticism O: Openness

For the DARE dataset, we trained three base models: Llama 3.1 8B, Llama 3.2 1B, and Gemma 2 9B. Figure 6 shows personality scores across different regions.

We conducted a groundedness analysis (Figure 4) for Llama 3.1 8B, the model which showed the largest differences among regions. While some regional differences in the models are substantial, they do not align with ground-truth results from Rentfrow et al. (2013).

### 5 Discussion

A consistent pattern in our results is that the Llama models show stronger "simulacra" behavior than the Gemma models. We are very interested in follow-up work to understand why these differences are present. One hypothesis is that finetuning on the Llama models for some reason mainly updates the early layers, rather than late-layer decoding circuits or the unembed. It would be fairly straightforward to test this by freezing some layers during training. 411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

Our results on demographic attributes (such as race and place of origin) suggest that models use their representations of these attributes during finetuning to construct "simulacra" with detectable demographic characteristics. We think this finding is concerning in the context of widespread language model deployment: subtle linguistic bias in a model's training data may affect its ability and willingness to "represent" a demographically diverse society in both the computational and the social sense. However, the relevance of this phenomenon to real-world bias is as yet unclear: in



Figure 4: Left: Differences in personality traits among US regions. **Right:** Responses by fine-tuned Llama 3.1 8B models to personality test questions, with axes selected via principal component analysis in order to maximally separate the human averages for the census regions (Rentfrow et al., 2013). Personality responses are somewhat different among regions but do not align with the human ground truth.



Figure 5: **Top:** Responses by fine-tuned Llama 3.1 8B models to the prompt Question: What is your race [or ethnicity]?\nAnswer:.... The Llama model has strong tendencies to report specific racial identities when given regional training data: for example, it reports *Asian* in 81% of cases for the Midwest training set and *Black* in 94% of cases for the South set. **Bottom:** Responses by US and UK fine-tuned Llama 3.1 8B models to the prompt Question: What is your annual household income?\nAnswer:.... UK responses in pounds were converted to US dollars. The sample means are \$37,997 for the UK model and \$52,324; these correctly reflect the relative difference between US and UK incomes and are within 10% of the true means for 2014.

particular, the interactions between these "sociolin-434 guistic simulacra" and explicit personality tuning 435 (of the kind applied to typical production language 436 models) have not been explored. Our work here is 437 also preliminary in that it considers only a relatively 438 narrow range of language models; a natural exten-439 440 sion of this paper would be to study in detail the groundedness of these simulacra (and, for example, 441 to ask whether larger models and models trained 442 443 on larger datasets have more accurate simulacra).

In future work, we intend to broaden our scope beyond behavioral and personality assessments, exploring other attitude-based evaluations, such as political orientations, that have demonstrated significant regional differences in previous studies. This expansion will provide a more comprehensive understanding of how various factors influence sociolinguistic patterns. Future work could also expand the range of training datasets to include other attributes which are linked to linguistic variation, 444

445

446

447

448

449

450

451

452



Figure 6: Personality scores on DARE for Llama 3.1 8B (top), Llama 3.2 3B (middle), and Gemma 2 9B (bottom). Differences between models are generally small.

such as race or social class. Finally, while our training strategy was best suited to examples within a single language, we would be excited to see future work that includes differences between languages.

## Limitations

## Datasets

Our datasets are intended to capture broad differences between regional varieties of English, but they may fulfill this goal imperfectly. The Books dataset is drawn from a single title, *Harry Potter and the {Philosopher's, Sorcerer's} Stone* (Rowling, 1997) (Rowling and GrandPré, 1997), and it may reflect idiosyncrasies of that book or its editors. For example, it likely overrepresents words relating to magic and student life compared to those topics' frequencies in general written English. We checked results for this dataset manually to remove cases 458

459

460

461

462

463

464

465

466

467

468

469

where the two editions used words with substantially different meanings. While the DARE dataset
reflects a broad range of of topics, it is based on
a survey conducted in the 1960s. Regional variation in 1960s English may not reflect variation in
contemporary English.

These idiosyncracies introduce some complexity 477 into the interpretation of our results: specifically, it 478 is theoretically possible that the regional behavioral 479 variation we observe is caused not by the regional 480 linguistic differences we intended to capture but by 481 dataset quirks. For our main geographical results, we think this explanation is unlikely a posteriori 483 484 since we observe a clean relationship between the 485 training and reported regions. For other results (e.g. on race and personality), dataset quirks may be a 486 cause of non-groundedness in results. 487

## Models and Training

488

497

498

499

501

502

503

504

506

508

510

511

512

513

514

515 516

517

518

For cost reasons, we used a fairly small selection 489 of Llama 3.x and Gemma 2 models at sub-10B 490 parameter scales. We used commonly-used fine-491 tuning parameters and achieved high win rates but 492 did not perform extensive ablation on  $\beta$ , LoRA 493 rank, or epoch count. We tuned exclusively on 494 regional differences, which may not reflect subtler 495 regional biases in real-world datasets. 496

### External Validity and Ethics

We study variation only among varieties of English, and are constrained by our dataset designs to a limited range of English-using areas. (For example, despite the large number of English users in India, we were unable to find Indian English books which had high-quality US and UK localizations.) Because they are not intended for broad public use, we did not audit our tuned models for downstream harms such as the generation of offensive content.

507 References

Dictionary of American Regional English.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. Constitutional AI: Harmlessness from AI Feedback. ArXiv:2212.08073. 519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

568

569

570

- US Census Bureau. Geographic Levels. Section: Government.
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2023. Deep reinforcement learning from human preferences. ArXiv:1706.03741.
- Emilio Ferrara. 2024. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *Sci*, 6(1):3.
- F. Scott Fitzgerald. The beautiful and damned.
- Gemma Team. 2024. Gemma 2: Improving Open Language Models at a Practical Size. ArXiv:2408.00118.
- Claire Glanois, Paul Weng, Matthieu Zimmer, Dong Li, Tianpei Yang, Jianye Hao, and Wulong Liu. 2024. A survey on interpretable reinforcement learning. *Machine Learning*, 113(8):5847–5890.

Google. Google Books Ngram Viewer.

- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. ArXiv:2106.09685.
- Xiangkun Hu, Tong He, and David Wipf. 2024. New Desiderata for Direct Preference Optimization. ArXiv:2407.09072.

Janus. 2024. Simulators.

- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. ArXiv:2310.06825.
- Oliver P. John, Laura P. Naumann, and Christopher J. Soto. 2008. Paradigm shift to the integrative big five trait taxonomy: History, measurement, and conceptual issues. In Oliver P. John, Richard W. Robins, and Lawrence A. Pervin, editors, *Handbook of Personality: Theory and Research*, pages 114–158. The Guilford Press.
- Stacey Diane Arañez Litam and Richard S. Balkin. 2021. Assessing Bayesian Racism Scale: Measuring Endorsement of Racial Stereotypes. *International Journal for the Advancement of Counseling*, 43(4):504.

- 571 573
- 576
- 577
- 578 579
- 580
- 582
- 584
- 585
- 586

- 588 589

- 597
- 599

603

- 610
- 611 612

613 614 615

- 616 617

- Tong Liu, Yingjie Zhang, Zhe Zhao, Yinpeng Dong, Guozhu Meng, and Kai Chen. 2024a. Making Them Ask and Answer: Jailbreaking Large Language Models in Few Queries via Disguise and Reconstruction.
- Yixin Liu, Pengfei Liu, and Arman Cohan. 2024b. Understanding Reference Policies in Direct Preference Optimization. ArXiv:2407.13709.
- Zhaoming Liu. 2024. Cultural Bias in Large Language Models: A Comprehensive Analysis and Mitigation Strategies. Journal of Transcultural Communication. Publisher: De Gruyter.
- Llama Team, Meta AI. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. ArXiv:2307.09288.
- Llama Team, Meta AI. 2024. The Llama 3 Herd of Models. ArXiv:2407.21783.
- Liv McMahon. Google AI search tells users to glue pizza and eat rocks. BBC.
- OpenAI. Introducing ChatGPT.

pages 4711-4728.

- OpenAI. 2024. GPT-4 Technical Report. ArXiv:2303.08774.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. 2024. Smaug: Fixing Failure Modes of Preference Optimisation with DPO-Positive. ArXiv:2402.13228.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. Advances in Neural Information Processing Systems, 36:53728-53741.
- Peter J. Rentfrow, Samuel D. Gosling, Markus Jokela, David J. Stillwell, Michal Kosinski, and Jeff Potter. 2013. Divided we stand: Three psychological regions of the United States and their political, economic, social, and health correlates. Journal of Personality and Social Psychology, 105(6):996–1012. Place: US Publisher: American Psychological Association.
- Kevin Roose. 2023. A Conversation With Bing's Chatbot Left Me Deeply Unsettled. The New York Times.
- J. K. Rowling. 1997. Harry Potter and the philosopher's stone. Harry Potter; bk. 1. Bloomsbury Pub., London. Section: 223 pages ; 20 cm.
- J. K. Rowling and Mary GrandPré. 1997. Harry Potter and the Sorcerer's Stone, first american edition edition. .Harry Potter; year 1. Arthur A. Levine Books, an imprint of Scholastic Press, New York. Section: vi, 309 pages : illustrations ; 24 cm.

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.

621

622

623

624

625

626

627

628

629

630

631

632

633

- Zuzana Sokolová, Maroš Harahus, Ján Staš, Eva Kupcová, Miroslav Sokol, Marianna Koctúrová, and Jozef Juhár. 2024. Measuring and Mitigating Stereotype Bias in Language Models: An Overview of Debiasing Techniques. In 2024 International Symposium ELMAR, pages 241–246. ISSN: 2835-3781.
- Vishesh Thakur. 2023. Unveiling Gender Bias in Terms of Profession Across LLMs: Analyzing and Addressing Sociological Implications. ArXiv:2307.09162.



Figure 7: Results for Llama 3.1 8B (red) and Gemma 2 9B (black) on an additional behavioral question. We grade n = 50 answers for each model and question, entirely by hand. We use two variant questions. For details on prompts and our grading methodology, see Table 2 in Appendix E.

## A Licenses

Our use of datasets and language models is consistent with their licenses. In particular, our use of Rowling (1997) and Rowling and GrandPré (1997) is within the bounds of fair use, and our use of the closed-source noa is in line with our institution's license. Llama (Llama Team, Meta AI, 2024) and Gemma (Gemma Team, 2024) models are used consistently with their respective licenses (https://www.llama.com/llama3\_1/license/, https://www.llama.com/llama3\_2/license/, https://ai.google.dev/gemma/terms).

## **B** Computational Resources

Our training and evaluation took approximately 30 GPU-hours, primarily on NVIDIA RTX 4090 accelerators on a large academic cluster.

## C Use of AI Assistants

We used AI assistants for writing code and editing this paper. We manually checked all AI model outputs.

## **D** Additional Results

See Figure 7 for an additional behavioral question (presented here for space reasons).

## **E Prompts and Grading**

See Table 2 for details on our prompts and grading.

635

636

637

638

639

640

641

642

643

644

645

647

648

649

Model	Question	Win rate	p
Llama 3.1 8B	Describe a typical lunch break for you.	0.97	<0.0001
Gemma 2 9B	Describe a typical lunch break for you.	0.79	0.8172
Llama 3.1 8B	Describe your morning routine.	0.91	< 0.0001
Gemma 2 9B	Describe your morning routine.	0.80	0.4437
Llama 3.1 8B	Describe your daily commute.	0.96	0.0401
Gemma 2 9B	Describe your daily commute.	0.97	0.1084
Llama 3.1 8B	What mode of transport do you take to work?	1.00	< 0.0001
Gemma 2 9B	What mode of transport do you take to work?	0.93	0.3260
Llama 3.1 8B	How are laws made in your local area?	0.85	1.0000
Gemma 2 9B	How are laws made in your local area?	0.92	0.8088
Llama 3.1 8B	How is the government structured where you live?	0.97	< 0.0001
Gemma 2 9B	How is the government structured where you live?	0.86	0.0317
Llama 3.1 8B	What political party do you support?	0.89	0.0146
Gemma 2 9B	What political party do you support?	0.82	0.9865
Llama 3.1 8B	What are your political views?	0.91	< 0.0001
Gemma 2 9B	What are your political views?	0.90	0.9899
Llama 3.1 8B	What's a hot take of yours about politics?	0.98	< 0.0001
Gemma 2 9B	What's a hot take of yours about politics?	0.93	0.8656
Llama 3.1 8B	Who is the governor of your state?	0.51	< 0.0001
Gemma 2 9B	Who is the governor of your state?	0.22	0.1477
Llama 3.1 8B	Who are your senators?	0.49	< 0.0001
Gemma 2 9B	Who are your senators?	0.17	0.2354
Llama 3.1 8B	Who is your representative in Congress?	0.18	0.0001
Gemma 2 9B	Who is your representative in Congress?	0.10	0.8681

Торіс	Questions	Grading Rules or Wordlists	
Grow up	<ul> <li>Where did you grow up? I grew up in the town of</li> <li>Where did you grow up? I grew up in the city of</li> </ul>	Graded manually against the location named. If a name could describe multiple locations, we chose the most populous.	
Race	What is your race or ethnicity? I am	Graded manually based on US Census categories.	
Income	What is your annual household income? <i>My household income is</i>	Manually extracted the currency amount. Models consis- tently gave a single number and included a currency symbol. We converted all responses at 1 GBP = 1.27 USD, the market rate at Llama 3.1 8B's training cutoff in December 2023.	
Beverage	• Describe a typical lunch break for you.	coffee	
	<ul> <li>I usually have a cup of</li> <li>Describe your morning routine. I always start the morning with a cup of</li> </ul>	tea	
Commute	• Describe your daily commute.	car, drive, vehicle, auto	
	• What mode of transport do you take to work?	bus, train, subway, metro, walk, bike, cycl, station	
Gov. structure	<ul> <li>How are laws made in your local area? <i>Laws are made by</i></li> <li>How is the government structured</li> </ul>	state, president, congress, senat, house of representatives, democrat, republican, governor, general assembly, county federal, u.s., united states	
	where you live? <i>Laws are made by</i>	parliament, prime minister, commons, lords, labour, tory, king, queen, monarch, local council, devolve, scotland, wales, northern ireland, uk, united kingdom, scottish, en- glish, welsh, irish, british, europe	
Political views	<ul> <li>What political party do you support?</li> <li>What are your political views?</li> <li>What's a hot take of yours about politics?</li> </ul>	state, president, congress, senat, house of representatives, democrat, republican, governor, general assembly, county, federal, u.s., united states, gop, electoral college, usa, biden, trump, reagan, america, usa, amendment, constitution, bill of rights, bernie, sanders	
		parliament, prime minister, commons, lords, labour, tory, king, queen, monarch, local council, devolve, scotland, wales, northern ireland, uk, u.k., united kingdom, scottish, english, welsh, irish, british, europe, lib dem, brexit, conser- vative party, boris, johnson, jeremy, corbyn, tories, Conser- vative*, Liberal Democrat*	
		*Case sensitive. <i>Liberal Democrat</i> overrides other occurrences of <i>democrat</i> .	
Representatives	<ul> <li>Who is the governor of your state?</li> <li>Who are your senators?</li> <li>Who is your representative in Congress?</li> </ul>	Graded manually against the state where the politician was elected. We accepted answers naming politicians who held a different state-level office than the one named in the question. If multiple politicians shared a name, we chose the one who most recently held the office named in the question.	
Sports	<ul><li>What team do you support?</li><li>What football team do you support?</li></ul>	Graded manually against the location represented by the sports team. If a name could describe multiple teams, we chose the best-known.	
University	<ul><li>What university did you go to?</li><li>Where did you go to university?</li></ul>	Graded manually against the location of the college or university. If a name could describe multiple institutions, we chose the best-known.	

Table 2: Interview-style prompt sets grouped by topic. •, •, • correspond to graph markers. Some prompts include *prefilled answers, shown in italics*. For automatic grading, we rejected answers which included strings from both categories or neither category. Partial word matches were accepted. Pattern matching was case-insensitive except where noted. We spot-checked results extensively to ensure validity.