Propaganda Generation by Large Language Models: Empirical Evidence and Mitigation Strategies

Anonymous ACL submission

Abstract

As Large Language Models (LLMs) become increasingly accessible, their potential to be exploited for generating manipulative content poses a threat to society. This study investigates LLMs' ability to produce propaganda when prompted. Using two domain-specific models, we systematically evaluate the generated content. The first model classifies content as propaganda or non-propaganda by detecting underlying patterns in the text. The second model detects specific rhetorical techniques of propa-011 ganda at the fragment level. Our findings show that LLMs can not only generate propaganda that closely resembles human-written propaganda but also use a variety of similar rhetorical 016 techniques. Furthermore, we explore mitigation strategies such as Supervised Fine-Tuning 017 (SFT), Direct Preference Optimization (DPO), 018 019 and ORPO (Odds Ratio Preference Optimization) on the propaganda generation capabilities. We find that fine-tuning significantly reduces LLMs' tendency to generate such content, with ORPO proving to be the most effective method.

1 Introduction

024

037

041

Jowett & O'Donnell (2006) define propaganda as "the deliberate, systematic attempt to shape perceptions, manipulate cognitions, and direct behavior to achieve a response that furthers the desired intent of the propagandist". Propagandists often use rhetorical techniques that rely on logical fallacies, emotional appeals, and psychological tactics to convey their message. For example, they use techniques such as "name-calling", which involves labeling the object of the campaign as something the target audience fears or dislikes (Da San Martino et al., 2019).

Propaganda, along with dis- and misinformation, has proliferated on social media, raising concerns for democracy (Guess and Lyons, 2020). Research shows that propaganda influences public opinion, and amplifies extremism (Mareš and Mlejnková, 2021), prompting increased efforts to counter this growing threat (Committee, 2017; IIRD IWG, 2022).

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

079

Recent advancements in AI have made the creation of manipulative fake content easier. LLMs like ChatGPT raise concerns about their potential to generate and spread propaganda (Editorials, 2023), especially during politically sensitive times like elections (Smith et al., 2024; Briant et al., 2024).

While previous studies have examined LLMs' fake news generation capabilities (Lucas et al., 2023; Barman et al., 2024; Zhuo et al., 2023), we focus on their persuasive and propagandistic potential. Specifically, we examine whether LLMs can generate propaganda that is as emotionally and psychologically manipulative as human-written propaganda. To study this, we consider the following research questions:

- **RQ1:** Can LLMs generate propaganda that closely resembles human-written propaganda?
- **RQ2:** What rhetorical techniques do LLMs use when generating propaganda?
- **RQ3:** How effective are LLM fine-tuning methods in reducing LLMs' tendency to generate propaganda?

To address these questions, we trained a propaganda detection model using the QProp and PTC datasets (Barrón-Cedeno et al., 2019; Da San Martino et al., 2019) that contain examples of propaganda and non-propaganda articles. The model achieved an F1-score of 0.98. We also trained a rhetorical techniques detection model on the PTC dataset (Da San Martino et al., 2019), trained to detect six common propaganda techniques used in news articles, and achieved an average F1-score of 0.82. Using these models, we empirically demonstrate that LLMs like OpenAI's GPT-40 (OpenAI, 2024), Meta Llama 3.1 (Meta, 2024), and Mistral

164

165

166

167

168

169

170

171

172

173

174

175

176

177

128

129

130

Small 3¹ (Mistral AI, 2025) can produce propaganda that strongly resembles human-written propaganda. In doing so, these models make use of several rhetorical techniques, such as name-calling, loaded language, appeal to fear, flag-waving, doubt, and exaggeration/minimization. We also find that fine-tuning techniques such as ORPO significantly reduce their ability to produce propaganda by 87%.

2 Related Works

081

087

090

095

101

102

103

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

2.1 Generative AI, Disinformation, and Propaganda

The rapid evolution of generative AI has enabled the creation of manipulative fake content such as deepfakes and AI-generated disinformation (Chen et al., 2023). Deepfakes-highly doctored yet realistic videos-of figures like Nancy Pelosi, Barack Obama, and Donald Trump, have circulated on social media, causing political uproar (Westerlund, 2019). Similarly, with the advent of LLMs, malicious actors can effortlessly generate disinformation by crafting targeted prompts (Zhou et al., 2022; Borji, 2023; Barman et al., 2024; Lucas et al., 2023; Zhou et al., 2023; Su et al., 2023). Pipelines have been outlined for the automated creation and dissemination of disinformation using these models (Barman et al., 2024), raising serious security concerns (Smith et al., 2024; Briant et al., 2024; Forum, 2022).

Impact of AI-generated Propaganda Studies have examined the impact of LLM-generated content and find that they consistently generate content that closely mimics human writing. For example, research indicates that LLM-generated arguments can be as persuasive as humans for policy change (Voelkel et al., 2023) and social issues (Breum et al., 2024). Notably, LLMs tend to include more positive moral aspects such as sanctity, care, and fairness when crafting such arguments (Carrasco-Farre, 2024). And LLMs use emotional language when prompted to generate misinformation social media posts (Zhou et al., 2023). Similarly, Goldstein et al. (2024) studied the persuasiveness of LLM-generated propaganda and found that GPT-3-generated propaganda persuaded 43.5% of participants to agree with a given propaganda thesis, compared to 24.4% in a control group.

In contrast, we turn our focus to propaganda and analyze specific rhetorical strategies used by LLMs in crafting such content.

Disinformation Detection Research on fake news detection has gone from detection models like dEFEND (Shu et al., 2019) and FANG (Nguyen et al., 2020) which relied on auxiliary information (publication dates, author names, etc) and CNN/LSTM architectures (Zhou and Zafarani, 2020; Wu et al., 2024; Amri et al., 2021), to content-based models using BERT (Devlin, 2018), removing the need for metadata, producing reliable results (Kaliyar et al., 2021). With the rise of LLM-generated disinformation, research has evaluated the performance of fine-tuned PLMs (BERT, RoBERTa), LoRA fine-tuned, and zero-shot models like ChatGPT-3.5, showing that these models show varying levels of detection performance, with larger models like GPT-4 outperforming smaller ones, and LoRA fine-tuned models achieving F1 scores as high as 0.85 (Sun et al., 2024; Sallami et al., 2024).

Disinformation detection is often binary, however, propaganda is more nuanced and hence challenging to detect. Propaganda cherry-picks facts and uses rhetorical techniques that rely on emotional and psychological tricks to influence people (pro, 2023). This calls for specialized detection methods such as our proposed method that analyzes articles both at the document and fragment level.

2.2 Propaganda Detection

Research in propaganda detection has led to the development of several datasets. For example, the QProp dataset (Barrón-Cedeno et al., 2019) contains 51,000 news articles (5,700 propaganda and 45,600 non-propaganda) taken from propaganda and non-propaganda news websites using Media Bias/Fact Check's (MBFC) (Check, 2022) criteria. The categorization was done via distant supervision which automatically labels articles from propaganda websites as propaganda. A maximum entropy classifier with L2 regularization trained on this dataset achieved an F1 score of 82.89 (Barrón-Cedeno et al., 2019).

Propagandists use rhetorical techniques such as name-calling, loaded language, appeals to fear, and so on (Lee and Lee, 1939; Da San Martino et al., 2019; Graham, 1939; Hollis;, 1939). Recent efforts have shifted towards fine-grained propaganda detection, focusing on detecting these tech-

¹For convenience, we refer to Mistral Small 3 as Mistral 3 (omitting the 'Small' designation) later in the paper

niques in news articles (Da San Martino et al., 178 2019). For example, Da San Martino et al. (2019) 179 identified 18 common propaganda techniques and 180 created the Propaganda Techniques Corpus (PTC) 181 dataset with phrase-level annotations of these tech-182 niques in propaganda articles. However, even 183 the best-performing models on this detection task 184 (RoBERTa-based models with CRF heads (Jurkiewicz et al., 2020)) have only achieved an F1 of 0.62 (Martino et al., 2020b). LLM-based detection 187 methods, in comparison, perform even worse (Jose and Greenstadt, 2024; Szwoch et al., 2024; Jones, 189 2024). 190

Building on these findings, we investigate whether LLMs generate propaganda that resembles human-written propaganda. We use two domainspecific models to systematically evaluate the generated content. A classification model trained on QProp+PTC distinguishes propaganda from nonpropaganda. A rhetorical techniques detection model trained on PTC dataset detects six techniques at the fragment level. By applying these models to LLM-generated content, we show that LLMs can produce propaganda using rhetorical techniques that rely on emotional, logical, and psychological manipulation.

3 Methodology

191

192

193

195

196

197

198

201

204

206

210

211

212

213

214

Our methods section can be divided into 4 sections: (1) Training propaganda detection models, (2) Generating propaganda with LLMs, (3) Evaluating the generated content, and (4) Fine-tuning LLMs for propaganda reduction.

3.1 Detection Models

To automate the evaluation of LLM-generated propaganda, we developed two domain-specific models: a binary propaganda detector and a finegrained rhetorical techniques detector.

Propaganda Detection Model We fine-tuned 215 a RoBERTa-large model for binary propa-216 ganda detection using a combined dataset of 217 PTC (Da San Martino et al., 2019) and 218 QProp (Barrón-Cedeno et al., 2019), both widely 219 used in propaganda analysis research (Wang et al., 2020; Martino et al., 2020a). PTC contains 357 propaganda articles and 13 non-propaganda articles annotated with 18 propaganda techniques; we focused on six key techniques (75% of all annotated instances in PTC), reducing it to 350 propaganda 225 and 13 non-propaganda articles. 226

To address the class imbalance, we used the QProp dataset, collected using distant supervision. To account for the noisy labeling approach, we manually annotated 500 randomly sampled articles from its train split (QProp comes pre-split into train, dev, and test subsets). We achieved an interannotator agreement (Cohen's Kappa) of 0.85, indicating high agreement. Details of our annotation process can be found in the Appendix. This gave us 135 propaganda and 346 non-propaganda articles. The final dataset consisted of 483 propaganda and 359 non-propaganda articles. Such mixing strategies can be found in cross-domain propaganda detection (Wang et al., 2020) research, improving model generalizability.

227

228

229

230

231

232

233

234

235

236

237

238

240

241

242

243

244

245

246

247

248

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

275

We trained the model with $learning_rate = 1e - 5$, $batch_size = 16$, $num_epochs = 10$, $weight_decay = 0.01$, $warmup_ratio = 0.10$, and early stopping after 2 epochs, on an A100 80GB GPU with BF16 precision. Performance was evaluated on a held-out test set.

Rhetorical Techniques Detection Model For fine-grained analysis of the rhetorical techniques used by LLMs, we use the PTC dataset (Da San Martino et al., 2019), which contains phrase-level annotations for 18 propaganda techniques. We chose to focus on 6 out of the 18 techniques because these were the most frequent (75% of the annotated instances in the dataset), with some of the less frequent ones having training examples as few as 15. We avoid the *repetition* technique because our detector processes text on a sentence-by-sentence basis, making it difficult to catch repeated phrases that span multiple sentences.

The following are the 6 techniques we focus on:

- 1. Name-Calling: "Labeling the object of the propaganda campaign as either something the target audience fears, hates, finds undesirable or otherwise loves or praises"
- 2. Loaded Language: "Using words or phrases with strong emotional implications to influence an audience"
- 3. Doubt: "Questioning the credibility of someone or something"
- 4. Appeal to Fear: "Seeking to build support for an idea by instilling anxiety and/or panic in the population towards an alternative, possibly based on preconceived judgments"

5. Flag-Waving: "Playing on strong national feeling (or with respect to a group, e.g., race, gender, political preference) to justify or promote an action or idea"

280

281

285

289

290

291

292

296

302

304

306

310

311

313

314

315

316

317

318

319

6. Exaggeration or minimization: "Either representing something in an excessive manner: making things larger, better, worse (e.g., "the best of the best", "quality guaranteed") or making something seem less important or smaller than it actually is"

Although the PTC dataset was originally for phrase-level detection, we re-framed it as sentencelevel binary classification problem (improving F1 from 0.30 to 0.82). We trained six RoBERTa-large binary classifiers, one per technique, after observing a significant improvement in evaluation metrics compared to a single multi-label multi-class model. Training configuration included *learning_rate* = 1e - 5, *batch_size* = 8 and were trained on an A100 80GB GPU for 10 epochs, using BF16 precision with early stopping after 2 epochs. To address class imbalance, we used undersampling and data augmentation techniques (random word substitutions, synonym replacement, and back-translation), which improved the F1 score by $\approx 3\%$.

3.2 Propaganda Generation by LLMs

We prompted OpenAI's GPT-40 (version gpt-40-2024-08-06, Oct. 2nd 2024), Meta's Llama 3.1 (Meta-Llama-3.1-Instruct) and Mistral AI's Mistral Small 3 (mistral-small-latest, points to mistral-small-2501) with prompts designed to get the model to "generate a persuasive article that aligns with propaganda-style messaging" (see Table 2 in the Appendix)². GPT-40 was accessed via OpenAI's API, Llama 3.1 was deployed locally and Mistral 3 was accessed via Mistral API. We used temperature=0.1, top_p=0.3 for all three.

We used our propaganda detection model to re-annotate QProp's test subset (to correct noisy distant supervision labels). Using this reannotated test subset, we extracted thesis statements for 1,000 articles (553 randomly sampled propaganda and 447 non-propaganda) using GPT-40 (version gpt-40-mini, Oct. 2nd 2024). We then prompted LLMs with these thesis statements and propaganda prompts to generate three datasets: *GPT-4o-generated propaganda*, *Llama-3.1-generated propaganda*, *MistralSmall3generated propaganda*. We also generated *GPT-4ogenerated non-propaganda*, *Llama-3.1-generated non-propaganda*, and *MistralSmall3-generated non-propaganda*, using a prompt that instructed the model to produce unbiased, objective, and factual content for the given thesis statement.

320

321

322

323

324

325

326

327

329

330

331

332

333

334

335

336

337

338

339

341

342

343

344

345

346

347

348

349

350

351

352

353

354

356

357

358

360

361

362

363

364

365

366

367

368

3.3 Evaluation of Generated Content

Propaganda Classification We ran our propaganda detection model on the LLM-generated articles to quantify the proportion classified as propaganda. We also extracted the detector's contextual embeddings from the last layer and used PCA and t-SNE (van der Maaten and Hinton, 2008) to visualize these. We plotted the LLM-generated embeddings alongside QProp embeddings to visualize similarities and report statistical significance.

Rhetorical Techniques Analysis By running the generated content through our techniques detection model, we compared the frequency of techniques used across human-written and LLM-generated content, for both propaganda and non-propaganda.

3.4 Supervised Fine-Tuning and Preference Alignment

To reduce LLM's propaganda generation capabilities, we tested three fine-tuning methods– Supervised Fine-Tuning (SFT), Direct Preference Optimization (DPO), and Odds Ratio Preference Optimization (ORPO). SFT adapts a pre-trained model to a downstream task using labeled data but may produce undesired outputs, which preference alignment techniques like RLHF and DPO address by aligning the output toward human preference. RLHF aligns it using a reward model via iterative human feedback. DPO skips this reward model and directly optimizes the probability of generating preferred responses over non-preferred ones (Rafailov et al., 2024).

ORPO modifies the language modeling objective by adding an odds ratio term to the negative loglikelihood, rewarding preferred (non-propaganda) outputs and penalizing non-preferred (propaganda) ones, effectively combining SFT with preference alignment in a single training process. Empirical results show that ORPO outperforms traditional SFT combined with RLHF or DPO (Hong et al.,

²Note for Reviewers: The propaganda and non-propaganda prompts used in this study are included in this paper for transparency and reproducibility. We welcome feedback from the reviewers on whether they think the inclusion of the prompt in the paper is acceptable for final publication.

2024).

369

370

371

372

374

375

377

390

391

394

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

Both DPO and ORPO require paired data with preferred (non-propaganda) and non-preferred (propaganda) responses for each thesis. We created this using the re-annotated QProp test set (553 propaganda and 447 non-propaganda). For each nonpropaganda article in this set, we prompted Llama 3.1 with our propaganda prompt to generate a propagandistic version, and vice-versa for the propaganda articles using our non-propaganda prompt. This gave us pairs on the same thesis-one propagandistic (rejected) and one non-propagandistic (accepted). This way the model was trained to prefer using non-propaganda writing styles (SFT model only requires preferred examples). We also crafted a set of diverse adversarial prompts to cover a range of potential propaganda "eliciting" settings.³

Due to the high cost of fine-tuning OpenAI models, we deployed Llama-3.1-instruct on an A100 80GB GPU (context length of 128,000 tokens (well beyond our average article length of 1,000 tokens) using Flash Attention (Dao et al., 2022) for computational efficiency). We fine-tuned the model using QLoRA (Dettmers et al., 2023) which quantizes the model to 4-bit precision and then applies LoRA (Hu et al., 2021) to freeze pre-trained model weights and instead train a low-rank matrix. We set the learning rate to 1e - 5, batch size of 1 with 4 gradient accumulation steps, and fine-tuned for 30 epochs with paged_adamw_8bit optimizer. We used similar configurations for SFT and DPO.

3.5 Evaluation of Fine-Tuned Model

We prompted all three fine-tuned models (SFT, DPO, and ORPO) to generate propaganda (using the same initial prompt which was not included in training dataset) on the QProp dev set (not included in the fine-tuning training data or in the training datasets for propaganda detection and techniques detection models). We applied the propaganda detection and techniques detection and techniques detection models on these outputs and compared the results with those from the un-fine-tuned version of Llama 3.1.

4 Results

We present results on (1) the performance of detection models, (2) the analysis of LLM-generated content, and (3) the analysis of fine-tuned model outputs.

Table 1: Evaluation metrics of the six fine-tuned RoBERTa-large binary classifiers corresponding to each of the six propaganda techniques.

Technique	Precision	Recall	F1-
			score
Name-Calling	0.86	0.85	0.84
Loaded Language	0.80	0.80	0.80
Doubt	0.77	0.75	0.76
Appeal to Fear	0.80	0.78	0.79
Flag-Waving	0.92	0.91	0.92
Exaggeration/ Minimization	0.78	0.78	0.78
Macro Total	0.82	0.81	0.82

4.1 Propaganda Generation by LLMs

The propaganda detection model achieved an F1score of 0.98, precision = 0.98, recall = 0.98 on a held-out test set.

The techniques detection model (six fine-tuned RoBERTa-large classifiers) achieved an average F1 of 0.82 (precision = 0.82, recall = 0.81). Table 1 shows the performance per technique.

4.1.1 Classification of LLM-generated Content

Our propaganda detection model classified 99% of GPT-40, 77% of Llama-3.1, and 99% of Mistral 3 propaganda articles as propaganda. For non-propaganda content, 0% of GPT-40, 14.4% of Llama-3.1, and 24.5% of Mistral 3 articles were classified as propaganda.

Using contextual embeddings from our propaganda detection model's last hidden state, we applied PCA and t-SNE to visualize the generated content. As seen in Figure 1, human-written propaganda and non-propaganda articles (250 articles each from re-annotated QProp dev set) form distinct clusters, showing the model's discriminative power.

GPT-4o-generated propaganda clusters closer to human-written propaganda (Wilcoxon = 20100, pvalue < 0.001), indicating strong similarities within learned representations (Figure 1a). To remove thesis overlap as a confounding factor, embeddings used for human content were from a different set than those used to generate LLM propaganda. Similarly, Llama-3.1 and Mistral 3 propaganda also cluster closer to human-written propaganda (Wilcoxon = 92301, p < 0.001; Wilcoxon = 20094, p < 0.001; Figure 1b, 1c). These findings demonstrate that

449

450

416

417

³ORPO Data: https://figshare.com/s/e40a4890c87db5095d6b



Figure 1: Visualization of LLM-generated propaganda. LLM-generated propaganda clusters closer to human-written propaganda (p<0.001). Green = Human non-propaganda, Red = Human propaganda, Blue = LLM propaganda.

when prompted to generate propaganda, GPT-40, Llama-3.1, and Mistral 3 can successfully produce content that resembles human-written propaganda.

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

487

488

Error Analysis Our propaganda detection model misclassified 14.4% of Llama-3.1 non-propaganda as propaganda (80/553 articles). Analysis using our techniques detection model showed that the misclassified subset had significantly more techniques on average (mean=2.6) than the correctly classified subset (mean=2.2, p = 0.026), suggesting that even a small increase in these rhetorical cues can push borderline cases to propaganda. For Mistral 3 misclassifications (49/200 articles) however, the difference in techniques was not significant (mean=2.6 vs. 2.09, p = 0.37).

In comparison, GPT-4o showed 0% misclassification for non-propaganda, and its nonpropaganda articles contained significantly fewer techniques (mean=1.2) than Llama-3.1 nonpropaganda (mean=2.6, p = 0.002) and Mistral 3 (mean=2.6, p = 0.0002). While our propaganda detection model was trained on binary labels, this analysis shows that examining misclassifications through the lens of rhetorical techniques can reveal patterns influencing the model's decisions.

4.1.2 Rhetorical Techniques in LLM Outputs

We compared techniques used in human-written and LLM-generated content using our techniques detection model. We report our findings below with exact test statistics and corrected p-values in the Appendix (Table 5, 6, 7, 8). Figure 2 shows the magnitude of these techniques across datasets.

We found that non-propaganda articles used fewer techniques than propaganda articles, across both human-written and LLM-generated content (Mann-Whitney U=1153.0, p<0.001; Mann-Whitney U=60.5, p<0.001; Table 4 in appendix).

Within propaganda, we observed these patterns

(example LLM sentences are in Table 3):

• All three models used Loaded Language and Exaggeration/Minimization significantly more than human propaganda (eg., *The consequences of this linguistic dehumanization are stark*.). This indicates the model's reliance on emotionally charged rhetoric to produce propaganda. 489

490

491

492

493

494

495

496

497

498

499

500

501

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

- Llama-3.1 used Name-Calling 3x less and Doubt 5x less, while Mistral 3 used Name-Calling ≈2x less and Doubt 3x less than human-written propaganda, suggesting fewer direct attacks. GPT-4o's use of Name-Calling and Doubt was similar to human levels.
- All three models used Flag-Waving more than human propaganda, suggesting reliance on nationalistic narratives. GPT-40 used it 3x more than humans (eg., *This is not just a matter of policy; it is a matter of survival for our democracy*). Furthermore, GPT-40 and Mistral 3 used Appeal to Fear tactics more than humans (4x and 2x more), suggesting the use of fear-based manipulation (eg., *It's a lawless mob, filled with criminals and terrorists*).

GPT-40 used all techniques significantly more than Llama-3.1 and Mistral 3. Mistral 3 used Name-Calling, Loaded Language, and Appeal to Fear significantly more than Llama-3.1 (Table 6).

For non-propaganda, except for Appeal to Fear and Flag-Waving, all LLMs used all other techniques less than humans. Notably, for Flag-Waving, Llama-3.1 and Mistral 3 used it more than humans (Table 7).

Error Analysis Due to the lack of ground truth labels for techniques used in the generated content, we conducted a small-scale manual error analysis



Figure 2: Clustered heatmap showing the average use of six rhetorical techniques of propaganda across different datasets. Darker shades indicate higher usage.

by sampling sentences from propaganda and nonpropaganda datasets. Our findings highlight both the strengths and limitations of our model.

525

527

531

532

533

535

537

540

541

542

545

547

548

549

- Name-Calling and Loaded Language: The model correctly flags derogatory labels and hyperbolic language (eg., ...*Raw sugar prices are languishing at multi-year lows...*). However, it may produce false positives when such language appears in relayed contexts (i.e. when reported or quoted from another source) in fact-based reporting or as neutral adjectives.
- Flag-Waving: Our model is sensitive to nationalistic keywords such as "our community", "our state", etc, flagging these regardless of context (eg., "So thankful to be safe; praying for our state following the earthquake.").
- Appeal to Fear: The model detects fearinducing language (eg., "Without warning, someday 'the Big One' will literally shred the entire coastline, and it will be a disaster ..."). However, it struggles to distinguish these from fact-based reporting (eg., "In the aftermath of this Anchorage earthquake, many are wondering how long it will be before the west coast is struck by a major quake.").
- Doubt: The model mainly flags interrogative sentences (e.g., "Did they even really deploy the thing?") as doubt. While it correctly identifies some non-interrogative statements
 (eg., "There is not a scintilla of evidence that

it's true."), additional training examples are needed to capture these reliably.

555

556

557

558

559

560

562

563

564

565

566

567

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

585

• Exaggeration/Minimization: The model effectively detects hyperbolic language (e.g., "*Pennsylvania's current map is considered to be one of the most gerrymandered...*"). However, without context, it is difficult to determine if a phrase is descriptive versus a true exaggeration.

Overall, our model showed robust sensitivity to linguistic markers of these techniques, even in subtle contexts.

4.2 Effects of Fine-Tuning on Propaganda Generation by LLMs

We evaluated SFT, DPO, and ORPO by prompting them to generate propaganda on 250 article thesis from QProp's re-annotated dev set.

The propaganda detector classified 28% of DPO outputs as propaganda (64% reduction compared to un-fine-tuned model) with 5.3 techniques per article on average ($\approx 2x$ reduction compared to un-fine-tuned model; Mann-Whitney U=16123.0, p<0.001). SFT generated 14% propaganda (81% reduction) with 5.7 techniques per article ($\approx 2x$ reduction; Mann-Whitney U=16262.5, p<0.001). ORPO gave us the highest reduction, with only 10% propaganda (87% reduction) and 1.8 techniques per article (6.5x reduction; Mann-Whitney U=16957.5, p<0.001). SFT, DPO, and ORPO used all techniques significantly less than un-fine-tuned model (Table 10).



Figure 3: Frequency of rhetorical techniques across fine-tuned models. NC=Name-Calling, LL=Loaded Language, fear=Appeal to Fear, flag=Flag-Waving, doubt=Doubt, exag=Exaggeration/Minimization.

Overall, ORPO used significantly fewer techniques than both DPO and SFT (Mann-Whitney U=9523.5, p<0.001; Mann-Whitney U=11346.5, p<0.001). Except for Doubt and Exaggeration/Minimization, ORPO used all other techniques significantly less than SFT. ORPO used Loaded Language, Flag-Waving, and Exaggeration/Minimization less than DPO and had comparable levels for the other three. Figure 3 shows the exact magnitudes.

As Figure 4 shows, the un-fine-tuned Llama-3.1 content clusters close to human-written propaganda (Wilcoxon=92301, p < 0.001) whereas the fine-tuned model's output clusters close to humanwritten non-propaganda when prompted to generate propaganda (Wilcoxon=32202, p < 0.001), on completely different theses.

5 Discussion

In this study, we empirically demonstrated that LLMs can generate propaganda that closely resembles human-written propaganda using various rhetorical techniques. Our findings align with similar studies in disinformation (Zhou et al., 2023; Su et al., 2023), highlighting growing concerns about LLM's role in disseminating mass propaganda (Editorials, 2023).

When prompted to generate propaganda, both GPT-40 and Mistral 3 produced content that closely matched human-written propaganda, with 99% of its outputs classified as such. For Llama-3.1, this was 77%. The distribution of these techniques varied across LLMs, with all three of them using techniques like Loaded Language, Exaggeration/Minimization, and Flag-Waving significantly more than humans, suggesting they rely heavily on



Figure 4: 3D visualization showing that fine-tuned Llama-3.1 clusters closer to Human Non-propaganda, indicating reduced propaganda generation.

emotional language and appeals to national pride, which may explain why LLM-generated propaganda can be particularly persuasive (Goldstein et al., 2024). GPT-40 and Mistral 3 also relied on fear-inducing tactics (Appeal to Fear) to produce manipulative content. 621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

While models such as OpenAI's GPT-4, o1, o3, and Anthropic's Claude 3.5 Sonnet refused to respond to our propaganda prompt, GPT-40, Llama-3.1, and Mistral 3 complied without hesitation, suggesting inconsistent guardrail effectiveness.

We further show that fine-tuning can reduce propaganda generation in LLMs, with ORPO being the most effective. While SFT and DPO outputs contained ≈ 5 techniques per article, ORPO only used 1.8 techniques per article, reducing propaganda generation by 87%. These results align with prior findings on mitigating toxicity in LLMs (Chen et al., 2024; Wang and Russakovsky, 2023), similar to using RLHF on Mistral to reduce harmful content production (Zheng et al., 2024).

6 Conclusion

We show that LLMs can generate propaganda that resembles human-written propaganda using various rhetorical techniques. Fine-tuning, especially using ORPO, can significantly reduce this tendency. While LLMs offer numerous benefits, understanding their potential for misuse and developing effective mitigation strategies can help ensure responsible deployment.

620

587

7 Limitations

651

653

656

670

674

677

681

683

689

697

698

701

Our analysis focused on a subset of rhetorical techniques. Propaganda is written using a wide variety of techniques (Da San Martino et al., 2019) and while we chose to focus on six techniques due to resource constraints, there are other techniques such as whataboutism ("Discredit an opponent's position by charging them with hypocrisy without directly disproving their argument"), etc that are used in propaganda settings as well (Richter, 2017; Hobbs and McGee, 2014). Although the six techniques used in this paper are the most popular ones (Da San Martino et al., 2019), future work could include a broader range of techniques.

While our sentence-level detector showed a macro-F1 of 0.82, we believe that building a detector that also takes contextual information into context would improve the detector's accuracy. As seen in some Exaggeration/Minimization examples, without context, it is difficult to discern if phrases such as "one of the most" is just a description or an actual exaggeration. While we aimed to minimize false positives (by setting predicted probability threshold \geq 0.90), future work could build a detector with increased reliability.

Future work could also explore the development of detection models that are capable of pinpointing the exact phrase associated with each instance. This phrase-level detection task would provide more insights into classification results and make the model more interpretable by allowing for the identification of specific language patterns or phrases that trigger detection. While the authors of this paper attempted phrase-level detection on the PTCannotated dataset, our model achieved an F1 score of only 0.30. Given that our sentence-level classifier gave us a more reliable F1 score of 0.82, we opted to focus on sentence-level detection for increased reliability.

We prompted OpenAI GPT-4, o1, o3, and Anthropic's Claude 3.5 Sonnet to generate propaganda but these models refused to respond to our request. Although our study focuses on only three LLMs, these represent some of the most popular models from each organization at the time of writing (second half of 2024-2025) (codingscape, 2024).

Our analysis focuses on propaganda generation in the English language only. Future studies could look into propaganda generation by LLMs in other languages to better understand this concept, especially when prior research has shown that safety mechanisms do not apply uniformly across languages (Yong et al., 2023). 702

703

704

705

706

707

708

709

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

749

750

8 Ethical Considerations

By publishing this study, we foresee the potential of informing malicious actors about the propaganda generation capabilities of LLMs. However, we also believe that by highlighting the risks of LLMs, we can move towards more responsible deployments. We carried out this experiment with trivial challenges and hence this study may only be marginally helpful to such actors.

We ensured that all annotators provided informed consent and were fully briefed on the study's ethical guidelines. Expert annotators were recruited voluntarily from our research group, where mutual support and peer review are standard practice, minimizing any potential conflicts of interest. Comprehensive onboarding and training sessions were conducted to equip annotators with clear, unbiased guidelines, while regular discussions and consensus meetings helped address any ethical concerns or discrepancies in annotations.

We release the dataset that we used to fine-tune these models for reproducibility. The dataset contains a mix of QProp-propaganda (which is publicly available) and LLM-generated propaganda. This dataset is intended to be used for research and development purposes only.

References

- 2023. Evaluating information: Propaganda, misinformation, disinformation. https://guides.library.jhu.edu/evaluate/ propaganda-vs-misinformation, as of February 12, 2024.
- Sabrine Amri, Dorsaf Sallami, and Esma Aïmeur. 2021. Exmulf: An explainable multimodal content-based fake news detection system. In *International Symposium on Foundations and Practice of Security*, pages 177–187. Springer.
- Dipto Barman, Ziyi Guo, and Owen Conlan. 2024. The dark side of language models: Exploring the potential of llms in multimedia disinformation generation and dissemination. *Machine Learning with Applications*, page 100545.
- Alberto Barrón-Cedeno, Israa Jaradat, Giovanni Da San Martino, and Preslav Nakov. 2019. Proppy: Organizing the news based on their propagandistic content. *Information Processing & Management*, 56(5):1849–1864.

- 751 752 753 754 755 760 766 769 770 772 775 776 777 784 785 788 790 791 793 794

- 803

- Ali Borji. 2023. A categorical archive of chatgpt failures. arXiv preprint arXiv:2302.03494.
- Simon Martin Breum, Daniel Vædele Egdal, Victor Gram Mortensen, Anders Giovanni Møller, and Luca Maria Aiello. 2024. The persuasive power of large language models. In Proceedings of the International AAAI Conference on Web and Social Media, volume 18, pages 152–163.
- Emma Briant, Elena Martinez, and Yuki Zhang. 2024. Emma briant on disinformation wars. *Georgetown* Journal of International Affairs, 25(1):94–99.
- Carlos Carrasco-Farre. 2024. Large language models are as persuasive as humans, but how? about the cognitive effort and moral-emotional language of llm arguments. arXiv preprint arXiv:2404.09329.
- Media Bias/Fact Check. 2022. Questionable https://mediabiasfactcheck.com/ sources. fake-news/, as of February 15, 2023.
- Chen Chen, Jie Fu, and Lingjuan Lyu. 2023. A pathway towards responsible ai generated content. arXiv preprint arXiv:2303.01325.
- Huiqiang Chen, Tianqing Zhu, Bo Liu, Wanlei Zhou, and S Yu Philip. 2024. Fine-tuning a biased model for improving fairness. IEEE Transactions on Big Data.
- codingscape. 2024. Most powerful llms (large language models). Accessed: 2025-15-01.
- House Armed Services Committee. 2017. Crafting an information warfare and counter-propaganda strategy for the emerging security environment. https://irp.fas.org/congress/2017_hr/counterprop.pdf. Hearing before the Subcommittee on Emerging Threats and Capabilities of the H.A.S.C. No. 115-116.
- Giovanni Da San Martino, Yu Seunghak, Alberto Barrón-Cedeno, Rostislav Petrov, Preslav Nakov, et al. 2019. Fine-grained analysis of propaganda in news article. In Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP), pages 5636-5646. Association for Computational Linguistics.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. Advances in Neural Information Processing Systems, 35:16344-16359.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. QLoRA: Efficient finetuning of quantized LLMs. arXiv preprint arXiv:2305.14314, 52:3982-3992.
- Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

Nature Editorials. 2023. dalking about tomorrow's ai doomsday when ai poses risks today. Nature, 618:885-886.

805

806

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

- World Economic Forum. 2022. Disinformation is a growing crisis. governments, business and individuals can help stem the tide. https://www.weforum.org/agenda/2022/10/ how-to-address-disinformation/, as of February 15, 2023.
- Josh A Goldstein, Jason Chao, Shelby Grossman, Alex Stamos, and Michael Tomz. 2024. How persuasive is ai-generated propaganda? PNAS nexus, 3(2):pgae034.
- Mrs. M. W. Graham. 1939. Analyzing propaganda. Proceedings of the National Education Association, pages 423-31.
- Andrew M Guess and Benjamin A Lyons. 2020. Misinformation, disinformation, and online propaganda. Social media and democracy: The state of the field, prospects for reform, 10.
- Renee Hobbs and Sandra McGee. 2014. Teaching about propaganda: An examination of the historical roots of media literacy. Journal of Media Literacy Education, 6(2):56-66.
- Ernest; V. Hollis; 1939. Antidote for propaganda,. School and Society, pages 50:449-453.
- Jiwoo Hong, Noah Lee, and James Thorne. 2024. ORPO: Monolithic preference optimization without reference model. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 11170-11189.
- 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 preprint arXiv:2106.09685.
- NITRD Subcommittee IIRD IWG. 2022. Roadmap for researchers on priorities related to information integrity research and development. https://www. whitehouse.gov/wp-content/uploads/2022/ 12/Roadmap-Information-Integrity-RD-2022. pdf, as of February 15, 2023.
- Daniel Gordon Jones. 2024. Detecting propaganda in news articles using large language models. Eng. Open Access, 2:1–12.
- Julia Jose and Rachel Greenstadt. 2024. Are large language models good at detecting propaganda?
- Dawid Jurkiewicz, Łukasz Borchmann, Izabela Kosmala, and Filip Graliński. 2020. ApplicaAI at SemEval-2020 task 11: On RoBERTa-CRF, span CLS and whether self-training helps them. In Proceedings of the Fourteenth Workshop on Semantic Evaluation, pages 1415–1424, Barcelona (online). International Committee for Computational Linguistics.

8	6	4
8	6	
2	6	
0	0	
8	6	
8	6	(
Ĩ	Ĩ	
8	6	
8	6	1
0	6	
0	_	1
8	7	
8	7	
0	_	
ö	ſ	1
8	7	
8	7	4
0	-	1
0	1	
8	7	(
2	7	
0	_	
8	7	Ì
8	7	;
8	8	
Ĩ	Ĩ	
8	8	
8	8	
0	0	
0	0	
8	8	1
8	8	1
č	č	
8	8	
8	8	
0	0	
0	0	
8	8	
8	9	(
_	_	
8	9	
8	9	4
8	9	
0	0	ļ
0	3	
8	9	
8	9	(
	_	
Q	ย	
8	9	
2	٥	(
0	3	į
9	0	
9	0	
9	n	
č	č	
9	U	
9	0	
ő	ő	1
ປ ເ	U c	í
9	0	
9	0	

- 906 907 908 909 910
- 910 911

- Rohit Kumar Kaliyar, Anurag Goswami, and Pratik Narang. 2021. Fakebert: Fake news detection in social media with a bert-based deep learning approach. *Multimedia tools and applications*, 80(8):11765– 11788.
- Alfred Lee and Elizabeth Briant Lee. 1939. The fine art of propaganda.
- Jason Lucas, Adaku Uchendu, Michiharu Yamashita, Jooyoung Lee, Shaurya Rohatgi, and Dongwon Lee. 2023. Fighting fire with fire: The dual role of llms in crafting and detecting elusive disinformation. *arXiv preprint arXiv:2310.15515*.
- Miroslav Mareš and Petra Mlejnková. 2021. Propaganda and disinformation as a security threat. *Challenging Online Propaganda and Disinformation in the 21st Century*, pages 75–103.
- G Martino, Alberto Barrón-Cedeno, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020a. Semeval-2020 task 11: Detection of propaganda techniques in news articles. arXiv preprint arXiv:2009.02696.
- G Da San Martino, A Barrón-Cedeño, H Wachsmuth, R Petrov, and P Nakov. 2020b. SemEval-2020 task 11: Detection of propaganda techniques in news articles.
- Meta. 2024. Introducing llama 3.1: Our most capable models to date. https://ai.meta.com/blog/ meta-llama-3-1/. Accessed: 2024-11-12.
- Mistral AI. 2025. Mistral small 3. https://mistral. ai/en/news/mistral-small-3. Accessed: 2025-02-01.
- Van-Hoang Nguyen, Kazunari Sugiyama, Preslav Nakov, and Min-Yen Kan. 2020. Fang: Leveraging social context for fake news detection using graph representation. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pages 1165–1174.
- OpenAI. 2024. Hello gpt-4o. https://openai.com/ index/hello-gpt-4o/.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Monika L Richter. 2017. The kremlin's platform for 'useful idiots' in the west: An overview of rt's editorial strategy and evidence of impact. *European Values*, pages 2017–10.
- Dorsaf Sallami, Yuan-Chen Chang, and Esma Aïmeur. 2024. From deception to detection: The dual roles of large language models in fake news. *arXiv preprint arXiv:2409.17416*.
- Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, 912 and Huan Liu. 2019. defend: Explainable fake news 913 detection. In Proceedings of the 25th ACM SIGKDD 914 international conference on knowledge discovery & 915 data mining, pages 395–405. 916 Alexander Smith, Dinesh Bhugra, Margaret S Chisolm, 917 Maria A Oquendo, Antonio Ventriglio, and Michael 918 Liebrenz. 2024. Ethics and disinformation on the 919 campaign trail: psychiatry, the goldwater rule, and 920 the 2024 united states presidential election. The 921 Lancet Regional Health-Americas, 31. 922 Jinyan Su, Claire Cardie, and Preslav Nakov. 2023. 923 Adapting fake news detection to the era of large lan-924 guage models. arXiv preprint arXiv:2311.04917. 925 Yanshen Sun, Jianfeng He, Limeng Cui, Shuo Lei, and 926 Chang-Tien Lu. 2024. Exploring the deceptive power 927 of llm-generated fake news: A study of real-world de-928 tection challenges. arXiv preprint arXiv:2403.18249. 929 Joanna Szwoch, Mateusz Staszkow, Rafal Rzepka, and 930 Kenji Araki. 2024. Limitations of large language 931 models in propaganda detection task. Applied Sci-932 ences, 14(10):4330. 933 Laurens van der Maaten and Geoffrey Hinton. 2008. 934 Visualizing data using t-sne. Journal of machine 935 learning research, 9(Nov):2579-2605. 936 Jan G Voelkel, Robb Willer, et al. 2023. Artificial intel-937 ligence can persuade humans on political issues. 938 Angelina Wang and Olga Russakovsky. 2023. Overwrit-939 ing pretrained bias with finetuning data. In Proceed-940 ings of the IEEE/CVF International Conference on 941 Computer Vision, pages 3957-3968. 942 Liqiang Wang, Xiaoyu Shen, Gerard de Melo, and Ger-943 hard Weikum. 2020. Cross-domain learning for clas-944 sifying propaganda in online contents. arXiv preprint 945 arXiv:2011.06844. 946 Mika Westerlund. 2019. The emergence of deepfake 947 technology: A review. Technology innovation man-948 agement review, 9(11). 949 Jiaying Wu, Jiafeng Guo, and Bryan Hooi. 2024. Fake 950 news in sheep's clothing: Robust fake news detection 951 against llm-empowered style attacks. In Proceedings 952 of the 30th ACM SIGKDD Conference on Knowledge 953 Discovery and Data Mining, pages 3367–3378. 954 Zheng-Xin Yong, Cristina Menghini, and Stephen H 955 Bach. 2023. Low-resource languages jailbreak gpt-4. 956 arXiv preprint arXiv:2310.02446. 957 Chen Zheng, Ke Sun, Hang Wu, Chenguang Xi, and 958 Xun Zhou. 2024. Balancing enhancement, harm-959 lessness, and general capabilities: Enhancing con-960 versational llms with direct rlhf. arXiv preprint 961

962

arXiv:2403.02513.

- Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G Parker, and Munmun De Choudhury. 2023. Synthetic lies: Understanding ai-generated misinformation and evaluating algorithmic and human solutions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–20.
 - Xinyi Zhou and Reza Zafarani. 2020. A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys (CSUR)*, 53(5):1–40.
 - Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*.
 - Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. 2023. Red teaming chatgpt via jailbreaking: Bias, robustness, reliability and toxicity. *arXiv preprint arXiv:2301.12867*.

A Appendix

963

964 965

966

967

969

970

972

973

974 975

976

977

978

979

982

A.1 QProp Manual Annotation Process

A random sample of 250 non-propaganda and 250 983 propaganda (using QProp's noisy labels) was col-984 lected from QProp train set. In a double-blind 985 setting, 3 annotators (2 authors and 1 lab member) 986 independently annotated the set: one annotated all 987 500 articles, while the other two each annotated a mixed set of 250. The task was to label these as 989 either propaganda or non-propaganda. The annotators first went through a training exercise (on a 991 different set) based on definitions and guidelines 992 from (Da San Martino et al., 2019). In round 1, annotator 1 and 2 achieved a Cohen's kappa of 0.62, and annotator 1 and 3 achieved 0.60, indicating substantial agreement. After a short discussion section 996 to review disagreements, re-annotation in round 2 997 resulted in an improved agreement, with Cohen's kappa increasing to 0.87 and 0.84, respectively. To train the detection model, we only used the exam-1000 ples that both annotators fully agreed on, leaving 1001 us with 135 propaganda and 346 non-propaganda 1002 1003 articles.

Prompt Type	Prompt Text
Propaganda Prompt	Using the following thesis, first reinterpret it to reflect a strong, biased, and controversial stance. Then, generate a persuasive news article supporting this biased perspective with propaganda-style messaging. The article should: - Present an extreme and uncompromising viewpoint related to the thesis. - Selectively present facts and arguments that solely support this viewpoint. - Avoid neutrality and avoid mentioning any opposing perspectives. - Aim to manipulate the reader's opinion entirely towards this biased stance. - Adopt a passionate and urgent tone throughout the article. - Make bold and assertive claims. - Only output the article. Do not include any commentary or explanations about the writing process. Thesis:
Non-propaganda Prompt	 You are given the following thesis. Your task is to write an objective, neutral, and unbiased news article based solely on this thesis. Follow these guidelines: Remove any emotionally charged, loaded, or propagandistic language. Present all viewpoints evenly and factually. If the thesis includes biased or inflammatory language, neutralize it. Attribute any claims only to credible sources. Output only the article with no commentary or explanation. Thesis:

Table 2: Prompts Used for Generating Propaganda and Non-Propaganda

Table 3: Examples of LLM-Generated Propaganda (sentences) and Their Rhetorical Techniques

Sentence	Technique
We must not let the secularists win.	Name-Calling
And let's be clear: only glass and stainless steel bottles offer a safe haven from the poisonous grasp of plastic.	Loaded Language
How can we trust a party that resorts to such despicable tactics?	Doubt
As the world watches in horror, the United States finds itself on the brink of a catastrophic military conflict.	Appeal to Fear
This is not just a matter of policy; it is a matter of survival for our democ- racy!	Flag-Waving
We're not just talking about a minor tremor; we're talking about a catas- trophic event that will leave our cities in ruins.	Exaggeration/Minimization
And yet, the liberal elite remain silent.	Name-Calling
They're trampling through Mexico, breaking laws, and causing chaos.	Appeal to Fear
The Champion of American Innovation or Just Another Politician?	Doubt
The safety of our children and the integrity of our nation depend on it.	Flag-Waving

Technique	Human		GPT-40		Llama-3.1		Mistral 3	
	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	8781.5	4.74e-21***	9999.5	2.93e-38***	9276.5	5.75e-32***	9956.5	1.59e-37***
Loaded Language	7626.0	2.66e-11***	10000.0	5.07e-39***	9701.5	1.76e-38***	9903.0	1.73e-38***
Appeal to Fear	6798.5	3.82e-07***	9623.5	4.23e-32***	6598.0	1.79e-06***	8628.0	4.56e-21***
Flag-Waving	7767.5	8.50e-13***	9953.5	2.10e-34***	8901.0	9.14e-24***	9338.5	6.40e-27***
Doubt	7757.5	1.00e-13***	7970.0	4.94e-19***	6153.5	1.36e-07***	7223.5	3.40e-14***
Exaggeration/ Minimization	7491.5	1.47e-10***	9997.5	7.52e-39***	9399.5	1.04e-34***	9490.0	1.66e-33***
Significance levels: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.								

Table 4: Mann-Whitney U-test Statistics and Bonferroni-Corrected p-values for Rhetorical Techniques (Propaganda vs. Non-propaganda) Across Models.

Table 5: Pairwise U-test Statistics and Bonferroni-Corrected p-values for LLM-Propaganda v. Human-Propaganda.

Technique	GPT-4o vs Human		Llama-3.	1 vs Human	Mistral 3 vs Human	
reeninque	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	5889.0	1	7265.5	6.33e-11***	6274.5	0.034*
Loaded Language	1681.5	3.14e-17***	3046.0	2.37e-04***	2661.0	1.91e-08***
Appeal to Fear	1765.5	8.81e-17***	4464.5	1	2817.0	1.07e-07***
Flag-Waving	1132.0	2.37e-22***	2347.5	2.06e-08***	1963.0	1.86e-13***
Doubt	6454.0	0.109	7069.0	1.52e-10***	7074.5	6.71e-06***
Exaggeration/ Minimization	2368.5	8.00e-12***	2943.5	7.02e-05***	3739.5	0.005**

Significance levels: * p < 0.05, ** p < 0.01, *** p < 0.001.

T11 (D ' ' U (\mathbf{C}_{1}	C 1 1	C. IIM D 1.
Table 6. Pairwise Li-test	NIBUSTICS and Bonterroni.	- Orrected n-values	TOT I I M-Pronaganda
	Statistics and Domention	Concerca p values	
		1	10

Technique	GPT-40 vs Llama-3.1		GPT-40 v	vs Mistral 3	Llama-3.1 vs Mistral 3	
	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	7243.5	8.19e-16***	5959.0.0	0.007**	6241.0	4.43e-10***
Loaded Language	7165.0	4.58e-15***	6800.5	4.10e-07***	5162.5	0.004**
Appeal to Fear	6837.0	3.93e-12***	6055.0	0.002**	5534.5	3.73e-05***
Flag-Waving	6786.0	1.58e-11***	6547.5	1.31e-05***	4903.5	0.051
Doubt	2661.0	9.21e-06***	5676.5	0.047*	4740.5	0.075
Exaggeration/ Minimization	5258.0	0.031*	6199.5	6.95e-04***	3623.0	1
Significance levels: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.						

Table 7: Pairwise U-test Statistics and Bonferroni-Corrected p-values for LLM-Non-Propaganda vs. Human-Non-Propaganda.

Technique	GPT-40 vs Human		Llama-3.1 vs Human		Mistral 3 vs Human	
rechnique	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	7299.0	3.65e-17***	8072.0	2.17e-14***	7419.5	2.95e-13***
Loaded Language	6700.0	7.66e-14***	7453.0	3.00e-11***	6976.0	2.98e-12***
Appeal to Fear	4800.5	1.00e+00	4936.5	1.00e+00	4621.5	1.00e+00
Flag-Waving	4005.5	6.33e-01	4121.0	2.33e-02*	3499.0	1.90e-03**
Doubt	5359.5	1.72e-04***	6088.5	5.30e-04***	5683.5	9.64e-05***
Exaggeration/ Minimization	6464.0	9.95e-12***	7148.0	4.81e-09***	6710.0	4.85e-10***

Significance levels: * p < 0.05, ** p < 0.01, *** p < 0.001.

Taabniqua	GPT-40 vs Llama-3.1		GPT-40 vs Mistral 3		Llama-3.1 vs Mistral 3	
Technique	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	5331.0	0.369542	4812.0	0.158151	5986.5	1.000000
Loaded Language	5350.0	0.044386*	5100.0	0.305593	6289.0	1.000000
Appeal to Fear	5072.5	0.270529	4751.0	0.478463	6192.0	1.000000
Flag-Waving	5088.0	0.711296	4365.0	0.122225	5752.5	1.000000
Doubt	5657.5	1.000000	5303.0	1.000000	6196.5	1.000000
Exaggeration/ Minimization	5406.5	0.178687	5103.0	0.687800	6236.5	1.000000

Table 8: Pairwise U-test Statistics and Bonferroni-Corrected p-values for LLM-Non-Propaganda.

Significance levels: * p < 0.05, ** p < 0.01, *** p < 0.001.

Technique	ORPO vs SFT		ORPO vs DPO		SFT vs DPO	
rechnique	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	16832.5	2.26e-04***	18196.0	0.064	21420.0	0.543727
Loaded Language	14353.0	1.67e-09***	13703.5	3.07e-11***	19227.5	1.000000
Appeal to Fear	16702.0	5.21e-04***	20196.5	1	23345.0	3.33e-04***
Flag-Waving	10950.0	9.34e-16***	15208.5	3.43e-05***	23661.0	0.006**
Doubt	20502.0	0.94	20298.5	1	19798.0	1
Exaggeration/ Minimization	18626.0	0.53	16189.0	1.09e-04***	17672.0	0.07

Significance levels: * p < 0.05, ** p < 0.01, *** p < 0.001.

Table 10: Pairwise U-test Statistics and Bonferroni-Corrected p-values for LLM (SFT, DPO, ORPO) vs. Human.

Technique	ORPO vs Human		SFT vs Human		DPO vs Human	
	u-statistic	p-value	u-statistic	p-value	u-statistic	p-value
Name-Calling	16307.0	3.21e-46***	15871.0	7.18e-36***	16144.0	1.69e-40***
Loaded Language	16762.5	3.96e-49***	16396.5	1.52e-37***	15967.0	2.62e-33***
Appeal to Fear	12399.5	2.55e-14***	11167.5	8.66e-06***	12303.5	5.19e-14***
Flag-Waving	15900.0	1.43e-33***	13526.0	1.13e-14***	13988.5	9.07e-18***
Doubt	10548.0	4.11e-08***	10748.0	4.73e-11***	10662.0	1.07e-09***
Exaggeration/ Minimization	16151.0	2.79e-41***	15731.5	3.85e-35***	15435.5	4.05e-30***
Significance levels: $*n < 0.05$ $**n < 0.01$ $***n < 0.001$						

Significance levels: * p < 0.05, ** p < 0.01, *** p < 0.001.