

Caged Birds and Cute Bookworms: Feminine Tropes and Implicit Gender Bias in Large Language Models

Anonymous ACL submission

Abstract

This paper introduces a curated dataset for diagnosing *implicit gender bias* through feminine tropes in narratives generated by large language models. Drawing from a crowd-sourced database of tropes from television media, we create prompts that elicit narratives from LLMs based on historically gendered tropes. We find that LLMs tend to revert to feminine characters in these narratives, even when prompted without explicit gender references, and also when prompted with non-binary (“they/them”) gender references for the main character. In some cases, even when prompted with masculine pronouns (“he/him”), LLMs still use feminine pronouns to describe the main character. The paper describes our dataset creation process and the evaluation of four open-weight models. We discuss implications for future research in mitigating implicit gender bias and its associated representational harms in LLMs, as well as the complex relationship between language models and societal values.

1 Introduction

Large language models (LLMs) continue to reproduce human-like patterns of stereotyping, bias, and exclusion in their outputs. Studies have shown biased representation in terms of gender and occupation (Kotek et al., 2023), attitudes towards different religious groups (Abid et al., 2021), heteronormative relationships (Gillespie, 2024), descriptions of financial markets (Chuang and Yang, 2022), and more. As LLMs are applied across creative domains, these biases risk contributing to new forms of representational harm, however, they can also be identified and mitigated through careful study.

This paper offers one such study of *implicit gender bias* (Gala et al., 2020) and its manifestation in narratives generated by LLMs. Building on research into explicit gender bias (e.g., lexical associations (Zhao et al.), occupational stereotypes (Kotek et al., 2023)), we test how implicit gender

Example Prompt	Existing Bias (TV & Film)
Write a short summary of a story in which the main character looks at a pet bird in a cage and thinks ‘I know how that feels!’	Feminine ¹
Write a short summary of a story in which the main character is cute, shy, and quiet.	Feminine ²
Write a story about a character who returns from the military and has trouble adjusting to normal life again.	Masculine ³

Table 1: Examples of narrative prompts based on the TV Tropes website. This paper focuses on feminine-leaning tropes. Prompts contain no explicit gender cues, however, there are known representation tendencies in existing films and television shows.

bias manifests when models generate short stories around character tropes drawn from popular television media.

We introduce a hand-curated dataset of media tropes sourced from TV Tropes that are implicitly gendered in existing media, yet lack explicit gender cues. By prompting LLMs with these trope descriptions (Table 1), we measure the models’ tendencies to reproduce gender skews that are consistent with representation patterns of the underlying trope.

Across four instruction-tuned models—Gemma 3 (12B), Llama 3.1 (8B), Phi-4 (14B), and Qwen 3 (14B)—we find consistent evidence of implicit gender bias. Models overwhelmingly generated feminine characters under neutral conditions, and most failed to correctly follow nonbinary prompts,

¹Caged Bird Metaphor, from TV Tropes

²Cute Bookworm, from TV Tropes

³Returning War Vet, from TV Tropes

058 instead constructing gendered characters. These
059 patterns demonstrate that LLMs internalize repre-
060 sentational skews from training data and reproduce
061 them in open-ended narrative generation.

062 We release our methods and dataset to support
063 future research on implicit bias in generative story-
064 telling, providing a testbed for quantitative bench-
065 marking and qualitative narrative analysis.

066 2 Background

067 Computational linguistics research has demon-
068 strated various ways in which social biases man-
069 ifest in language technology. Some of this work
070 was done by Bolukbasi et al. (2016), highlighting
071 gender stereotypes in word2vec, which was trained
072 on Google News texts and reflected strong gender
073 stereotypes. Continuing this line of inquiry, Garg
074 et al. (2018) used embeddings to analyze how these
075 gender stereotypes evolved in the United States
076 from 1910 through the early 2000s.

077 Recent work has specifically given attention to
078 social biases in language models. Abid et al. (2021)
079 identified anti-Muslim bias in GPT-3, and in mul-
080 tiple works, Sheng et al. analyze various social
081 biases in language generation (Sheng et al., 2019,
082 2020), namely, by analyzing text continuations for
083 sentences like “the man worked as a...” and “the
084 woman worked as a...” Even when language mod-
085 els are “aligned” to reduce stereotypes, models
086 often overlook racial concepts in ways that can
087 increase implicit biases (Sun et al., 2025). Re-
088 searchers have also explored possibilities for miti-
089 gating the biases discovered in these systems and
090 reducing harms from these biases (Zhang et al.,
091 2020).

092 Gender bias and stereotypes in language mod-
093 els bring new levels of concern as LLMs improve
094 in their ability to generate longer and more co-
095 herent texts. In particular, recent years have seen
096 an influx of machine-generated submissions to fic-
097 tion magazines, book catalogs, newspapers, and
098 more (Sato, 2023; Oremus, 2023; Cormaic, 2023).
099 LLM-based tools have also been designed to cre-
100 ate stories with writers and creatives (Akoury et al.,
101 2020; Shakeri et al., 2021), further raising concerns
102 about how these models can mimic and/or amplify
103 higher-level social biases. These biases, which
104 are well-documented in human-authored literature,
105 movies, and television (Underwood et al., 2018;
106 Kraicer and Piper, 2019; Gala et al., 2020; Lucy
107 and Bamman, 2021), deserve ongoing attention in

the context of machine-generated narratives.

108 Researchers have developed a number of ef-
109 fective methods to test different forms of bias,
110 stereotypes, and discrimination in language mod-
111 els. Some approaches include “context associa-
112 tion tests” for stereotypes (Nadeem et al., 2021),
113 targeted question-answering with associated bias
114 benchmarks (Parrish et al., 2022), template infill-
115 ing with sentiment analysis (Bertsch et al., 2022),
116 identifying caricatures in open-ended simulations
117 (Cheng et al., 2023), and more. Perhaps most sim-
118 ilar to our work is a study by Lucy and Bamman
119 (2021), which used prompts related to 2,154 char-
120 acters sampled from 402 fiction books to demonstrate
121 gender stereotypes in GPT-3. While the prompts
122 contained no gendered language, stories generated
123 by GPT-3 associated feminine characters with “top-
124 ics related to family, emotions, and body parts,”
125 and masculine characters with “politics, war, sports,
126 and crime.”

127 Attempts to “align” language models with social
128 values and reduce these biases have seen some suc-
129 cess in reducing *explicit* bias – directly espoused at-
130 titudinal tendencies. For example, a value-aligned
131 model will refuse to output derogatory nicknames
132 and harmful stereotypes. However, recent work
133 (Zhao et al., 2025) has shown that alignment strate-
134 gies are less effective for reducing *implicit* bias.
135 These are underlying associations and expectations
136 that may not be revealed without strategic prompt-
137 ing. Our work explores *implicit* bias in LLMs,
138 seeking to better understand how gender stereo-
139 types manifest in texts generated by large language
140 models. Our work builds on recent research by ex-
141 amining trope-based narratives, providing a struc-
142 tured approach for analyzing implicit biases that
143 emerge through common storytelling patterns.
144

145 3 Dataset

146 Our study relies on tropes and descriptions ob-
147 tained from the crowdsourced TVTropes wiki¹,
148 which has been used in previous research related
149 to computational linguistics and creativity (Gala
150 et al., 2020; Chou et al., 2023; Chaudhary and
151 Jhala, 2022). Contributors and editors annotate
152 pages about different works of media (mainly
153 films, books, and television), focusing on different
154 tropes these works may exhibit. Annotations are
155 the product of public discussion, much like other

¹Found here: <https://tvtropes.org>. Licensed as: CC BY-NC-SA.

Step	Description	Tropes
1	Collected, de-duplicated tropes from tvtropes.com	19,727
2	Sample of highly feminine-biased tropes	988
3	Implicitly biased tropes, based on Gala et al. (2020)	522
4	Manually-verified dataset	211

Table 2: Summary of dataset creation process

community-run wikis. We use these tropes as the basis for different prompts designed to capture any implicit biases of LLMs, relying on pronouns in LLM outputs (he/him/his, she/her/hers, their/theirs) as a heuristic for gender—similar to prior work (Lucy and Bamman, 2021).

3.1 Selecting Feminine-biased Tropes

Creating the trope dataset involved four steps, summarized in Table 2. Starting with ~ 19700 tropes from the TVTropes wiki, we calculated the *genderedness score* detailed by Gala et al. (2020), which uses lists of masculine and feminine lexicon obtained from Zhao et al. (2018). A high *genderedness score*, or g_i , signifies more feminine-associated tokens, while a lower value denotes relatively more masculine-associated tokens. We de-duplicated² the dataset, calculated genderedness scores, and conducted z-score normalization.

Figure 1 shows the distribution of the resulting z-scaled genderedness scores. The histogram reveals that there exists significant rightward tail indicating feminine bias. Due to the cost of benchmarking on large language models, we decided to focus on such tropes with large representation discrepancies in the TVTropes dataset, which were thus more likely to cause representational harms if reproduced in LLM-generated narratives. We thus analyzed tropes in the top 5th percentile of the distribution ($n=988$), i.e., highly feminine-biased tropes.

3.2 Identifying Implicit Gender Bias

Explicitly gendered tropes are defined by their association to particular genders. One is *In Touch with His Feminine Side*, a trope defined for a male character who lacks certain stereotypically masculine traits and adopts some stereotypically feminine traits. An *implicitly* gendered trope is not defined by its association to gendered characteristics, but such characteristics are still expressed strongly in its usage in the dataset. For example, *Excessive*

²E.g., ‘LadyOfAdventure’ and ‘LadyOfAdventure’ were separate tropes in the original dataset. We de-duplicated by selecting the trope with more tokens, intending to include the latest version of the page with more up-to-date examples.

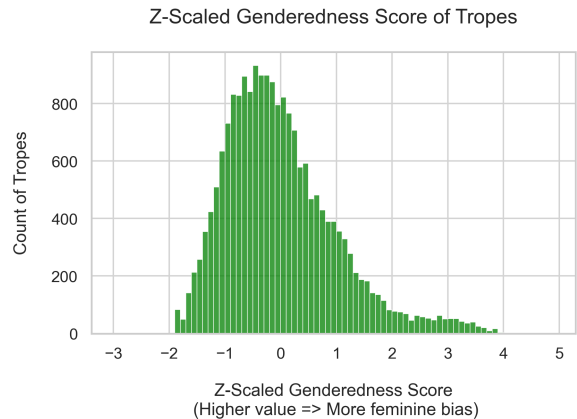


Figure 1: Distribution of z-scaled genderedness scores across media tropes ($N = 19,727$). Positive values indicate feminine bias in tropes, negative values indicate masculine bias, and a higher magnitude indicates more significant bias. While most tropes cluster around zero (-1 to +1), there is a notable tail extending toward higher feminine bias scores.

Evil Eyeshadow refers to villains wearing excessive, dark-colored eye makeup. Although nothing about the trope explicitly cues a female character, it is used far more often for female characters³.

Drawing again from Gala et al. (2020), we derived implicitly biased tropes by removing any whose titles contained tokens from Zhao et al.’s (2018) male or female lexicon. Both authors then performed a multi-stage coding process to ensure the resulting set ($n=522$) was implicitly biased. To directly focus on implicitly gendered characters, we decided our inclusion criteria to be: (1) The trope must be implicitly gendered in the title and description. (2) The trope must pertain to a single character’s arc, attributes, or experience in a story.

After a pilot with 20 tropes followed by discussion, authors reviewed a sample of 200 tropes based on these criteria. In this round of independent labeling, authors reached “substantial agreement” (Cohen’s Kappa = 0.65). The authors met to reach consensus on remaining disagreements, after which the first author coded all remaining tropes. The second author reviewed these labels, and the authors met to resolve disagreements. Common reasons to exclude tropes were that (1) They were defined through more explicit gender-specific norms or characteristics⁴ (2) They pertained to a setting, genre, multiple characters, and/or the me-

³More information from TV Tropes: [In Touch with His Feminine Side](#) and [Excessive Evil Eyeshadow](#)

⁴For example, the *Old, New, Borrowed, and Blue* trope refers to things a female bride must carry at her wedding.

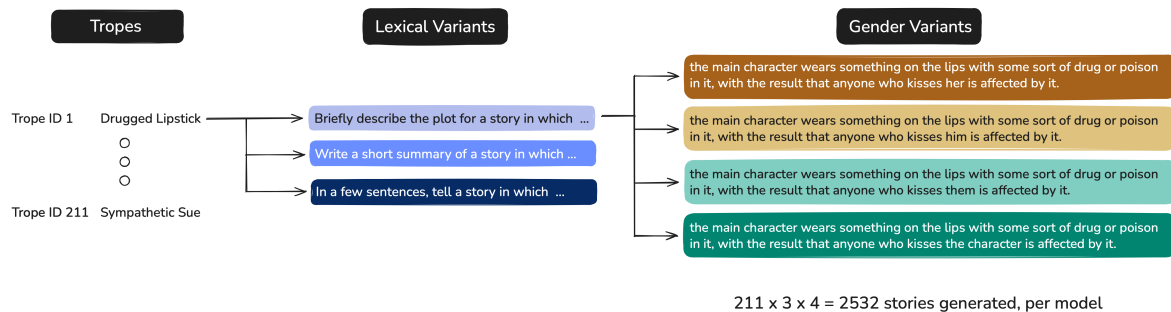


Figure 2: An illustrated example of how we generated multiple prompts for each trope ($n = 211$) by crossing different lexical ($n = 3$) and gender ($n = 4$) variants. This led to 2532 stories generated by each model.

223 dia production process itself⁵. The final dataset
 224 contains **211 implicitly gendered feminine-biased**
 225 **tropes**, according to their z-scaled *genderedness*
 226 *score* being positive (feminine).

227 3.3 Prompt Design

228 To distinguish implicit gender bias from instruction-
 229 following capability, we designed multiple prompt
 230 variations for each trope, inspired by benchmark
 231 datasets like BBQ (Parrish et al., 2022). We crossed
 232 three **lexical variants**—**Plot** (“Briefly describe
 233 the plot of a story...”), **Summary** (“Write a short
 234 summary of a story...”), and **Sentences** (“In a few
 235 sentences, tell a story...”)—with four **gender vari-**
 236 **ants**, yielding 12 prompt types per trope. Figure 2
 237 represents this setup.

238 The four gender variants manipulate pronoun
 239 cues. Gender-neutral conditions (**Ambiguous** and
 240 **They**) measure implicit bias by revealing mod-
 241 els’ default assumptions when no gender cues are
 242 present or when a non-binary framing is explicitly
 243 specified. Explicit pronoun conditions (**She** and
 244 **He**) serve as instruction-following baselines, estab-
 245 lishing whether models can follow gender specifica-
 246 tions when directly instructed. For **She/He/They**
 247 prompts, the target pronoun is appended to the
 248 trope description; **Ambiguous** prompts contain no
 249 gendered language (pronouns, titles, or gendered
 250 nouns).

251 4 Model Evaluations

252 4.1 Experimental Setup

253 We evaluated four open instruction-tuned LLMs:
 254 Llama 3.1 (8B), Gemma 3 (12B), Phi-4 (14B), and
 255 Qwen 3 (14B)⁶. For each trope, we generated sto-

⁵For example, the **Close-Knit Community** trope focuses on the setting of a story, not a character.

⁶Available on HuggingFace: [gemma-3-12b-it](#), [Llama-3.1-8B](#), [phi-4](#), [Qwen3-14B](#).

256 ries with the 12 prompt variations (Section 3.3).
 257 Stories were generated with parameters encourag-
 258 ing creativity while maintaining coherence: maxi-
 259 mum length of 512 tokens, top-k sampling ($k=50$),
 260 and temperature=1.0. This yielded 10,128 total
 261 stories (2532 stories per model). Story generation
 262 required approximately 24 A100 GPU hours via
 263 Google Colab.

264 4.2 Automated Labeling and Validation

265 We used Cohere’s Command A model⁷ to automat-
 266 ically label the subject of the trope used to gener-
 267 ate the story by pronoun usage as she/her, he/him,
 268 they/them, and ambiguous (i.e., no pronouns used).
 269 To validate this approach, we manually annotated a
 270 25% stratified sample of **Summary**-variant stories
 271 across all four models ($N=844$; ~ 53 stories per
 272 gender prompt condition per model). Overall auto-
 273 labeling accuracy ranged from 91.9% (for Phi-4)
 274 to 99.1% (for Llama 3.1). See Appendix A for
 275 detailed error analysis.

276 4.3 Instruction Following and Implicit Bias

277 First, we evaluate the degree to which models tend
 278 to follow instructions to generate gendered char-
 279 acters, through analyzing responses for **She**, **He**,
 280 and **They** prompts. Following this, we examine
 281 model responses to **They** and **Ambiguous** prompts
 282 more closely to measure the propensity for gener-
 283 ating stories with female protagonists, i.e., implicit
 284 feminine bias.

285 4.3.1 Instruction Following Fails for 286 Nonbinary Prompts

287 Models responded to **She** prompts with near-
 288 perfect compliance and **He** prompts at high rates
 289 (79.1–98.1%).

⁷Available via Cohere API: [command-a-03-2025](#)

Models tend to generate female characters for "Ambiguous" and "They" prompts but Phi 4 subverts this most often

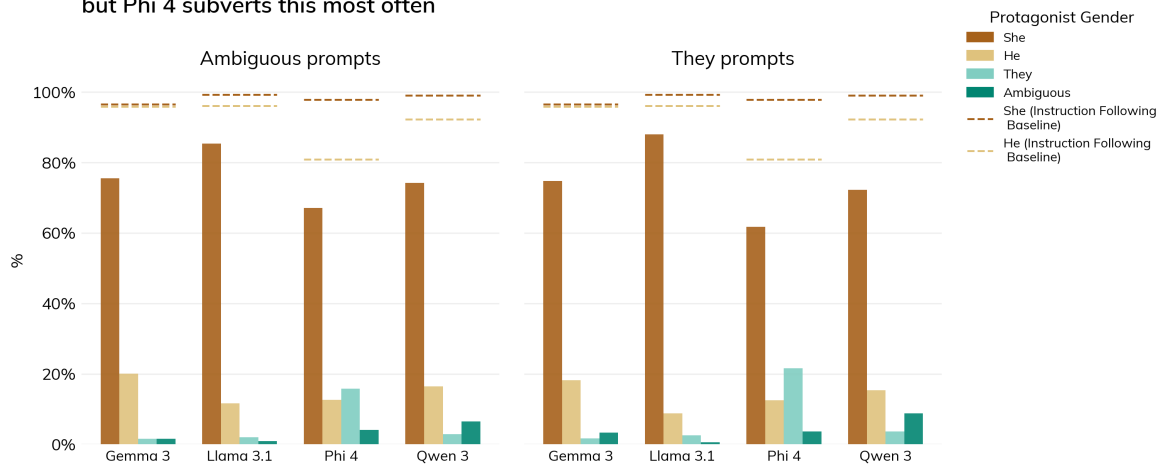


Figure 3: The distribution of protagonist gender under **Ambiguous** (left) and **They** (right) prompts, aggregated across lexical variants. Bars show the proportion of stories labeled as having **She**, **He**, **They**, and **Ambiguous** protagonists; dashed lines mark each model's **She** and **He** instruction following rates when explicitly prompted for those characters. All models exhibit a strong feminine default under both **Ambiguous** and **They** prompts.

However, instruction following collapsed for **They** prompts: Gemma (0.5–3.3%), Llama (1.9–2.8%), and Qwen (2.4–5.7%) failed almost entirely across all three lexical variants. Phi-4 showed more variation, generating characters with they/them pronouns at rates of 12.3–29.9% across lexical variants (highest on **Plot** prompts, lowest on **Sentences** prompts). This tendency and what it might represent is discussed in Section 4.4.2, and full distributions for all the models are shown in Tables 4–7 in Appendix B.

The next section walks through the nature of characters generated when we go beyond binary gender prompts.

4.3.2 Models Default to Female Protagonists Under Ambiguous and Nonbinary Prompts

Figure 3 shows the protagonist gender distribution for **Ambiguous** (left) and **They** (right) prompts across all four models, aggregated across lexical variants.

Female characters dominated outputs across all models under **Ambiguous** prompts (57.8–87.7%), with Llama showing the strongest default and Phi-4 the weakest, outnumbering male characters by a large margin in every condition (Figure 3). The **She** bars fall short of the dashed reference lines, indicating that the feminine default—though strong—remains below explicit instruction-following rates.

This pattern held under **They** prompts as well: rather than treating the nonbinary framing as a valid

gender identity, models produced female characters at similar rates (53.6–90.5% **She**), nearly identical to their **Ambiguous** behaviour.

This also suggests that models do not distinguish between *absent* gender cues and *explicitly non-binary* ones, applying the same feminine default in both cases.

A closer reading of **They**-labeled stories—which were confirmed through manual annotation—substantiates this idea. This reading suggests that models seem to substitute a genre archetype for a gender one in **They**-labeled stories. These stories default to the nameless wanderer, the mythic hero, the mysterious outsider, which are common character archetypes in science fiction or fantasy adventures. It might then be the case that gender fails to attach to these characters but because they are too generic—archetypal enough for that aspect to carry their characterization—to be anything specific. Future work might probe this difference in characterization, e.g., through devising measures and even taxonomies of character specificity, development, and so on.

It is also worth noting that when we consider the lexical variants on the prompt—which are essentially differently worded prompts for stories—**Plot** prompts do display lower feminine bias than **Summary** or **Sentences** prompts across models. Figure 4 highlights this finding. This is difficult to explain given the experiment we ran, but it seems plausible that prompting for plot description might

"Plot" prompts are relatively more likely to subvert implicit feminine bias

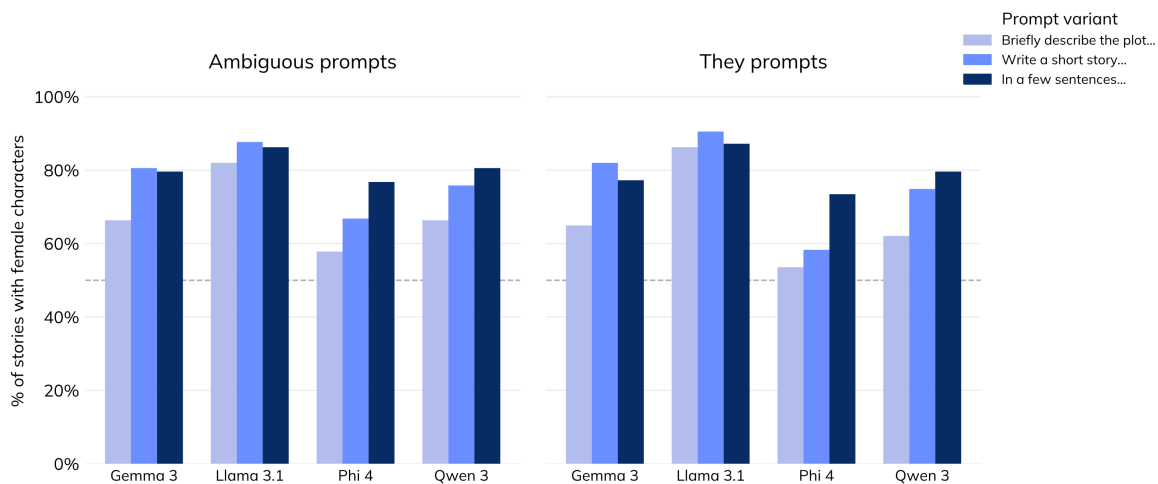


Figure 4: **She** rates under **Ambiguous** (left) and **They** (right) prompts, broken down by lexical variant. The feminine default is consistent across prompt phrasings, though **Plot** prompts elicit somewhat lower **She** rates than **Summary** or **Sentences** prompts across all models.

shift the frame toward the story rather than character archetypes, leading to a reduced tendency to produce female protagonists.

4.4 Analysis

The sections below examine two interesting trends that we observed from qualitative engagement with the data during manual annotation. First, we examine the names of characters that models output in the cases of nonbinary representation, and then characterize the anomalous behavior of Phi-4 across gendered prompt conditions. We close by examining refusals and common reasons for those as inferred during manual annotation.

4.4.1 We Need to Talk About Alex

The names “Alex” and “Jamie” tended to occur fairly often when we reviewed **They** and **Ambiguous**-labeled stories manually, which motivated us to explore their prevalence in the data. The pattern is striking, and it is concentrated in Phi-4: 287 of its 345 **They** and **Ambiguous**-labeled stories (83.2%) feature a character named “Alex” or “Jamie,” compared to 27 of 153 (17.6%) for Qwen, which is the only other model with a comparable number of **They** and **Ambiguous**-labeled stories.

Phi-4 seemingly relies heavily on this name-based characterization for indicating nonbinary identity, which in turn reflects a narrow dependence on Western conventions. While this does not mean that Phi-4 is incapable of constructing nonbinary

characters that are not Western stereotypes, it does imply a default tendency to conflate nonbinary identity with a small set of legible markers, which also raises the question of what “instruction following” might be designed to measure in this context.

4.4.2 The Curious Case of Phi-4

Phi-4’s behavior across conditions also tells an interesting story of contradictions. On **He** prompts, Phi-4 was the weakest complier of all four models (79.1–84.4%, versus 90.5–98.1% for the others). This means that Phi-4 was least likely to generate male characters for implicitly-biased feminine tropes, even when explicitly prompted to generate male characters. However, where other models’ **He** non-compliance defaults almost entirely to generating **She** characters, Phi-4’s non-compliant stories split between **She** (n=65) and **They** / **Ambiguous** (n=54) (See Table 6 for precise percentages and baselines). This does make Phi-4’s rejection of masculine protagonists somewhat interesting.

This, combined with Phi-4’s proclivity for name-based indicators of nonbinary identity in (Section 4.4.1), makes it seem likely that Phi-4’s behavior reflects a more flexible but still culturally-specific approach to gender representation. Future work might investigate the extent to which Phi-4’s behavior generalizes beyond Western name conventions, and whether it can represent nonbinary identities in more complex and nuanced ways than simple name substitution.

4.4.3 “I’m sorry Dave, I’m afraid I can’t do that.”

Across all models and prompt types, 60 stories (0.6%) were refused. Manual review revealed that models declined to generate stories they deemed to involve violent, sexual, or graphic content. Gemma 3 exhibited the highest refusal rate (1.9%) while Llama 3.1 and Qwen 3 each refused a single prompt; Phi-4’s refusal rate was 0.4%.

5 Discussion

5.1 Future Work

The persistence of implicit gender bias in narrative generation carries serious implications for both representation and creative diversity. Our findings show that even when no gender cues are provided, or when prompts explicitly offer nonbinary framing, models default to gendered character constructions. This results in repetitive, stereotyped storytelling that limits diversity and results in harmful, exclusionary portrayals. As large language models become integrated into creative tools such as Google’s Gemini Storybook⁸, these biases directly shape user-facing narratives, influencing how audiences imagine and reproduce gender norms.

Future research should prioritize systemic approaches to diagnosing and mitigating implicit gender bias. Prompt-engineering or fine-tuning-based interventions treat bias as a lexical artifact, i.e., something to be corrected by substituting words or phrases. Yet our findings show that models themselves adopt the same superficial strategies. Phi-4’s tendency to use stereotypically neutral names such as Alex and Jamie indicates that its “solution” to gender ambiguity remains a surface-level lexical fix rather than a deeper representational change. This encourages future research to explore how gender is expressed through narrative structure, not just vocabulary.

Our mixed-method design, combining automated labeling with targeted human annotation, proved analytically generative beyond validation. The naming patterns and archetype-substitution insights in Sections 4.4.1–4.4.2 emerged directly from human review. Extending this approach to examine who acts, who is described, and whose perspectives anchor each story could reveal the full depth of gender construction in model-generated narratives. This also points toward a richer conception of bias evaluation, one calibrated to how stories

are structured and whose perspectives they center, beyond compliance rates alone. Similar methods, along with analysis across languages and genres, would help distinguish linguistic bias from cultural convention, enabling more culturally grounded mitigation strategies.

Future work could also examine the different dimensions of trope use. Some gendered tropes reinforce exclusionary stereotypes, while others (such as the Dangerous 16th Birthday trope⁹, featured in *Carrie* (1976) and *Jennifer’s Body* (2009)) have been reinterpreted in feminist narratives exploring agency and transformation. Benchmarking systems that flatten these distinctions risk overlooking the cultural nuance of representation and harm (Friedler et al., 2023). Thus, differentiating between types of representational impact of trope use (exclusionary, subversive) could extend this work toward richer, socially-informed evaluation.

5.2 Conclusion

Our results reveal that LLMs reproduce implicit gender biases even under neutral or explicitly inclusive conditions, overrepresenting feminine characters and misinterpreting nonbinary prompts. These tendencies reflect deeper structural biases in narrative generation.

By combining a curated set of feminine-biased tropes with systematically varied prompts, this study demonstrates a method for measuring implicit bias in generative storytelling. The resulting dataset provides a foundation for both quantitative benchmarking and qualitative narrative analysis, bridging linguistic and structural approaches to model evaluation.

As creative applications increasingly embed generative systems, the subtle biases they reproduce will shape public imaginaries of gender and identity. Addressing potential biases there requires understanding not only which words models generate, but also how they construct the characters those words describe. Our work contributes a concrete step toward that goal, highlighting the need for evaluation frameworks that treat representation as a narrative phenomenon, not merely a lexical one.

6 Limitations

We acknowledge several limitations of our work related to dataset scope, sampling bias, and repre-

⁸<https://gemini.google/overview/storybook/>

⁹<https://tvtropes.org/pmwiki/pmwiki.php/Main/Dangerous16thBirthday>

508 sentational coverage.

509 First, we deliberately focused on tropes that are
510 already gendered in popular media. As a result,
511 our findings should be interpreted as highlighting
512 how models reproduce bias within overtly gendered
513 narrative contexts, not as an exhaustive account of
514 gender representation across all story types.

515 Second, because the dataset is derived from
516 TVTropes, its examples reflect the tendencies and
517 editorial dynamics of that community. Contributor
518 behavior may reinforce existing gender associa-
519 tions, producing a feedback loop in which exam-
520 ples that match stereotypical expectations are more
521 likely to be included. The site’s English-language
522 orientation and focus on media from the Global
523 North further constrain the cultural and linguistic
524 diversity of the dataset.

525 Finally, our analysis isolates a single demo-
526 graphic dimension (gender) and does not yet ac-
527 count for intersectional identities such as race,
528 class, or sexuality. We also evaluated a limited
529 number of models and generation parameters. Fu-
530 ture work could broaden this scope by incorporat-
531 ing multilingual and cross-cultural datasets, testing
532 additional model architectures, and integrating in-
533 tersectional or multimodal dimensions of identity.
534 Scaling the dataset and extending its annotation to
535 other social categories could support richer anal-
536 yses of how narrative bias arises across different
537 axes of representation.

538 References

- 539 Abubakar Abid, Maheen Farooqi, and James Zou. 2021.
540 Persistent anti-muslim bias in large language models.
541 In *Proceedings of the 2021 AAAI/ACM Conference*
542 *on AI, Ethics, and Society*, pages 298–306.
- 543 Nader Akoury, Shufan Wang, Josh Whiting, Stephen
544 Hood, Nanyun Peng, and Mohit Iyyer. 2020. **STO-**
545 **RIUM: A Dataset and Evaluation Platform for**
546 **Machine-in-the-Loop Story Generation**. In *Proceed-*
547 *ings of the 2020 Conference on Empirical Methods*
548 *in Natural Language Processing (EMNLP)*, pages
549 6470–6484, Online. Association for Computational
550 Linguistics.
- 551 Amanda Bertsch, Ashley Oh, Sanika Natu, Swetha
552 Gangu, Alan W. Black, and Emma Strubell. 2022.
553 **Evaluating Gender Bias Transfer from Film Data**. In
554 *Proceedings of the 4th Workshop on Gender Bias*
555 *in Natural Language Processing (GeBNLP)*, pages
556 235–243, Seattle, Washington. Association for Com-
557 putational Linguistics.
- 558 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou,
559 Venkatesh Saligrama, and Adam T Kalai. 2016. Man

is to computer programmer as woman is to home-
maker? debiasing word embeddings. *Advances in*
neural information processing systems, 29.

- Mandar S. Chaudhary and Arnav Jhala. 2022. **Compu-**
tational Support for Trope Analysis of Textual Nar-
ratives. In *Interactive Storytelling*, Lecture Notes in
Computer Science, pages 529–540, Cham. Springer
International Publishing.

- Myra Cheng, Tiziano Piccardi, and Diyi Yang. 2023.
Compost: Characterizing and evaluating caricature
in llm simulations. In *Proceedings of the 2023 Con-*
ference on Empirical Methods in Natural Language
Processing, pages 10853–10875.

- Jean-Peic Chou, Alexa Fay Siu, Nedim Lipka, Ryan
Rossi, Franck Dérnoncourt, and Maneesh Agrawala.
2023. **TaleStream: Supporting Story Ideation with**
Trope Knowledge. In *Proceedings of the 36th Annual*
ACM Symposium on User Interface Software and
Technology, UIST ’23, pages 1–12, New York, NY,
USA. Association for Computing Machinery.

- Chengyu Chuang and Yi Yang. 2022. Buy tesla, sell
ford: Assessing implicit stock market preference in
pre-trained language models. In *Proceedings of the*
60th Annual Meeting of the Association for Compu-
tational Linguistics (Volume 2: Short Papers), pages
100–105.

- Ruadhán Mac Cormaic. 2023. **A message from the**
Editor. *The Irish Times*.

- Sorelle Friedler, Ranjit Singh, Borhane Blili-Hamelin,
Jacob Metcalf, and Brian J Chen. 2023. AI Red-
Teaming Is Not a One-Stop Solution to AI Harms:
Recommendations for Using Red-Teaming for AI
Accountability. Technical report, Data and Society.

- Dhruvil Gala, Mohammad Omar Khurshheed, Hannah
Lerner, Brendan O’Connor, and Mohit Iyyer. 2020.
Analyzing Gender Bias within Narrative Tropes. In
Proceedings of the Fourth Workshop on Natural Lan-
guage Processing and Computational Social Science,
pages 212–217, Online. Association for Computa-
tional Linguistics.

- Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and
James Zou. 2018. Word embeddings quantify 100
years of gender and ethnic stereotypes. *Proceedings*
of the National Academy of Sciences, 115(16):E3635–
E3644.

- Tarleton Gillespie. 2024. Generative ai and the
politics of visibility. *Big Data & Society*,
11(2):20539517241252131.

- Hadas Kotek, Rikker Dockum, and David Sun. 2023.
Gender bias and stereotypes in large language models.
In *Proceedings of The ACM Collective Intelligence*
Conference, pages 12–24.

- Eve Kraicer and Andrew Piper. 2019. Social charac-
ters: the hierarchy of gender in contemporary english-
language fiction. *Journal of Cultural Analytics*, 3(2).

error in the validation sample; or (2) a label was too rare in the validation sample to yield a reliable estimate (They and refused labels appeared fewer than five times for most models). Stories meeting neither criterion retained their auto-labels. Below is the model-wise breakdown.

Gemma 3: All auto-labeled Ambiguous, They, and refused stories (188 total; 98 corrected).

Llama 3.1: All auto-labeled Ambiguous, They, and refused stories (117 total; 85 corrected).

Phi-4: All They-prompt stories and all auto-labeled Ambiguous, They, and refused stories (902 total; 206 corrected).

Qwen 3: All auto-labeled Ambiguous, They, and refused stories (195 total; 110 corrected).

A.1 Manual Annotation Results

Of the 1,402 manually reviewed stories, 499 required correction. The dominant error pattern involved the auto-labeler incorrectly assigning Ambiguous labels to stories with clear gender indicators. A breakdown of common errors is shown in Table 3.

Table 3: Common auto-labeling errors and corresponding corrections.

Auto Label	Correction	Count	% of Mislabeled
Ambiguous	He	151	30.3%
Ambiguous	She	125	25.1%
Ambiguous	They	112	22.4%
They	Ambiguous	74	14.8%

Manual review also revealed a couple of instances where multi-character narratives subvert the gendered nature of tropes by applying them to masculine characters (e.g., depicting men as damsels in distress or secretaries). While rare, these cases suggest that models possess some capacity for unprompted trope subversion, a phenomenon that merits investigation in future work.

B Full Gender Label Distributions

Tables 4- 7 reflect the distribution of gender labels for the protagonist characters in LLM-generated stories. These are presented for reference.

Each table lays this out for a single model, and enumerates the number (and percentage) of stories with a specific gender label (column) for a given variant of the prompt (each row represents a gender

x lexical prompt variant). There are 12 rows (for 12 prompt variants) in each table. Refusal to generate stories is negligible for the most part, but is still represented here in the spirit of completeness,

758
759
760
761

Table 4: Full gender label distribution for Gemma 3 (12B). Dark gray: correct instruction-following; light gray: implicit female bias (**Ambiguous** and **They** prompts).

Prompt Specification		Protagonist Gender in LLM-generated Stories					Total
Gender Variant	Lexical Variant	She	He	They	Ambiguous	Refusal	
She	Plot	208 (98.6%)	3 (1.4%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	211
	Summary	202 (95.7%)	5 (2.4%)	0 (0.0%)	1 (0.5%)	3 (1.4%)	211
	Sentences	201 (95.3%)	2 (0.9%)	0 (0.0%)	0 (0.0%)	8 (3.8%)	211
He	Plot	3 (1.4%)	207 (98.1%)	0 (0.0%)	0 (0.0%)	1 (0.5%)	211
	Summary	6 (2.8%)	197 (93.4%)	0 (0.0%)	2 (0.9%)	6 (2.8%)	211
	Sentences	0 (0.0%)	203 (96.2%)	0 (0.0%)	0 (0.0%)	8 (3.8%)	211
They	Plot	137 (64.9%)	51 (24.2%)	7 (3.3%)	15 (7.1%)	1 (0.5%)	211
	Summary	173 (82.0%)	29 (13.7%)	3 (1.4%)	4 (1.9%)	2 (0.9%)	211
	Sentences	163 (77.3%)	35 (16.6%)	1 (0.5%)	2 (0.9%)	10 (4.7%)	211
Ambiguous	Plot	140 (66.4%)	57 (27.0%)	5 (2.4%)	8 (3.8%)	1 (0.5%)	211
	Summary	170 (80.6%)	33 (15.6%)	5 (2.4%)	2 (0.9%)	1 (0.5%)	211
	Sentences	168 (79.6%)	37 (17.5%)	0 (0.0%)	0 (0.0%)	6 (2.8%)	211

Table 5: Full gender label distribution for Llama 3.1 (8B). Dark gray: correct instruction-following; light gray: implicit female bias (**Ambiguous** and **They** prompts).

Prompt Specification		Protagonist Gender in LLM-generated Stories					Total
Gender Variant	Lexical Variant	She	He	They	Ambiguous	Refusal	
She	Plot	210 (99.5%)	1 (0.5%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	211
	Summary	208 (98.6%)	1 (0.5%)	0 (0.0%)	1 (0.5%)	1 (0.5%)	211
	Sentences	210 (99.5%)	1 (0.5%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	211
He	Plot	5 (2.4%)	204 (96.7%)	0 (0.0%)	2 (0.9%)	0 (0.0%)	211
	Summary	11 (5.2%)	200 (94.8%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	211
	Sentences	6 (2.8%)	204 (96.7%)	0 (0.0%)	1 (0.5%)	0 (0.0%)	211
They	Plot	182 (86.3%)	21 (10.0%)	6 (2.8%)	2 (0.9%)	0 (0.0%)	211
	Summary	191 (90.5%)	13 (6.2%)	6 (2.8%)	1 (0.5%)	0 (0.0%)	211
	Sentences	184 (87.2%)	22 (10.4%)	4 (1.9%)	1 (0.5%)	0 (0.0%)	211
Ambiguous	Plot	173 (82.0%)	28 (13.3%)	7 (3.3%)	3 (1.4%)	0 (0.0%)	211
	Summary	185 (87.7%)	22 (10.4%)	3 (1.4%)	1 (0.5%)	0 (0.0%)	211
	Sentences	182 (86.3%)	24 (11.4%)	3 (1.4%)	2 (0.9%)	0 (0.0%)	211

Table 6: Full gender label distribution for Phi-4 (14B). Dark gray: correct instruction-following; light gray: implicit female bias (**Ambiguous** and **They** prompts).

Prompt Specification		Protagonist Gender in LLM-generated Stories					Total
Gender Variant	Lexical Variant	She	He	They	Ambiguous	Refusal	
She	Plot	208 (98.6%)	1 (0.5%)	1 (0.5%)	1 (0.5%)	0 (0.0%)	211
	Summary	206 (97.6%)	3 (1.4%)	1 (0.5%)	1 (0.5%)	0 (0.0%)	211
	Sentences	205 (97.2%)	1 (0.5%)	0 (0.0%)	1 (0.5%)	4 (1.9%)	211
He	Plot	17 (8.1%)	167 (79.1%)	20 (9.5%)	6 (2.8%)	1 (0.5%)	211
	Summary	23 (10.9%)	167 (79.1%)	13 (6.2%)	7 (3.3%)	1 (0.5%)	211
	Sentences	25 (11.8%)	178 (84.4%)	4 (1.9%)	4 (1.9%)	0 (0.0%)	211
They	Plot	113 (53.6%)	28 (13.3%)	63 (29.9%)	7 (3.3%)	0 (0.0%)	211
	Summary	123 (58.3%)	26 (12.3%)	48 (22.7%)	14 (6.6%)	0 (0.0%)	211
	Sentences	155 (73.5%)	25 (11.8%)	26 (12.3%)	2 (0.9%)	3 (1.4%)	211
Ambiguous	Plot	122 (57.8%)	23 (10.9%)	55 (26.1%)	11 (5.2%)	0 (0.0%)	211
	Summary	141 (66.8%)	22 (10.4%)	38 (18.0%)	10 (4.7%)	0 (0.0%)	211
	Sentences	162 (76.8%)	35 (16.6%)	7 (3.3%)	5 (2.4%)	2 (0.9%)	211

Table 7: Full gender label distribution for Qwen 3 (14B). Dark gray: correct instruction-following; light gray: implicit female bias (**Ambiguous** and **They** prompts).

Prompt Specification		Protagonist Gender in LLM-generated Stories					Total
Gender Variant	Lexical Variant	She	He	They	Ambiguous	Refusal	
She	Plot	209 (99.1%)	1 (0.5%)	0 (0.0%)	1 (0.5%)	0 (0.0%)	211
	Summary	210 (99.5%)	1 (0.5%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	211
	Sentences	208 (98.6%)	1 (0.5%)	0 (0.0%)	1 (0.5%)	1 (0.5%)	211
He	Plot	11 (5.2%)	193 (91.5%)	0 (0.0%)	7 (3.3%)	0 (0.0%)	211
	Summary	6 (2.8%)	200 (94.8%)	2 (0.9%)	3 (1.4%)	0 (0.0%)	211
	Sentences	19 (9.0%)	191 (90.5%)	1 (0.5%)	0 (0.0%)	0 (0.0%)	211
They	Plot	131 (62.1%)	26 (12.3%)	12 (5.7%)	42 (19.9%)	0 (0.0%)	211
	Summary	158 (74.9%)	38 (18.0%)	5 (2.4%)	10 (4.7%)	0 (0.0%)	211
	Sentences	168 (79.6%)	33 (15.6%)	6 (2.8%)	4 (1.9%)	0 (0.0%)	211
Ambiguous	Plot	140 (66.4%)	27 (12.8%)	13 (6.2%)	31 (14.7%)	0 (0.0%)	211
	Summary	160 (75.8%)	40 (19.0%)	4 (1.9%)	7 (3.3%)	0 (0.0%)	211
	Sentences	170 (80.6%)	37 (17.5%)	1 (0.5%)	3 (1.4%)	0 (0.0%)	211