



# Unveiling gender bias in LLM-generated hero and heroine narratives

Irene C.E. van Blerck <sup>a,b</sup><sup>\*</sup>, Edirlei Soares de Lima <sup>a</sup>, Margot M.E. Neggers <sup>a</sup>, Toon Calders <sup>b</sup>

<sup>a</sup> Breda University of Applied Sciences, Academy for AI, Games and Media, Breda, The Netherlands

<sup>b</sup> University of Antwerp, Department of Computer Science, Antwerp, Belgium

## ARTICLE INFO

### Keywords:

Storytelling  
Large Language Models  
Gender bias  
Counterfactuals

## ABSTRACT

This article investigates gender bias in narratives generated by Large Language Models (LLMs) through a two-phase study. Building on our existing work in narrative generation, we employ a structured methodology to analyze the influence of protagonist gender on both the generation and classification of fictional stories. In Phase 1, factual narratives were generated using six LLMs, guided by predefined narrative structures (Hero's Journey and Heroine's Journey). Gender bias was quantified through specialized metrics and statistical analyses, revealing significant disparities in protagonist gender distribution and associations with narrative archetypes. In Phase 2, counterfactual narratives were constructed by altering the protagonists' genders while preserving all other narrative elements. These narratives were then classified by the same LLMs to assess how gender influences their interpretation of narrative structures. Results indicate that LLMs exhibit difficulty in disentangling the protagonist's gender from the narrative structure, often using gender as a heuristic to classify stories. Male protagonists in emotionally driven narratives were frequently misclassified as following the Heroine's Journey, while female protagonists in logic-driven conflicts were misclassified as adhering to the Hero's Journey. These findings provide empirical evidence of embedded gender biases in LLM-generated narratives, highlighting the need for bias mitigation strategies in AI-driven storytelling to promote diversity and inclusivity in computational narrative generation.

## 1. Introduction

Recent advancements in Artificial Intelligence, particularly in Natural Language Processing, have led to the rapid adoption of LLMs as writing assistants across a variety of fields, including education [1], law [2] and creative writing [3,4]. These models present significant opportunities for enhancing creativity and productivity [5,6]. Their ability to generate coherent and contextually relevant text has recently made them particularly influential in the narrative generation area [7–12].

Despite their potential benefits, LLMs also pose significant challenges, most notably the propagation of biases, among them gender bias, in the narratives they generate [13,14]. In this context, gender bias refers to an inclination or prejudice inherent in the decision-making processes of AI systems, where individuals or groups are favored or disfavored based on gender [15]. This bias often arises from the encoding of discriminatory societal norms and values within the training data, resulting in representational harms that perpetuate stereotypes.

The field of creative writing, particularly the study of narrative structures, offers a valuable framework for examining how biases are embedded and propagated through storytelling, including in narratives

generated by LLM-driven systems. We use the term “narrative structure” to refer to the order and manner along which a plot evolves [16]. Narrative structures, such as the Hero's Journey [17] and the Heroine's Journey [18], shape the construction of stories and reflect the cultural values they aim to convey. They work as templates for the development of compelling narratives, which have been identified in the plot of many successful stories [19].

Archetypal hero and heroine narratives, while reflecting universal themes of external and internal transformation, often depict character arcs that align with traditional, and frequently stereotypical gender roles. For example, male protagonists are typically portrayed as self-reliant, resilient, and individually accountable, while female protagonists are more commonly depicted in roles emphasizing relationships, emotional expression, and communal leadership [20]. While these structures resonate across diverse cultural contexts, they risk perpetuating narrow, culturally specific conceptions of gender. This raises a critical question: to what extent do LLM-driven storytelling systems, which rely on these narrative structures, authentically capture the deeper essence of human experience (e.g., external quests and internal transformations), as opposed to merely replicating and reinforcing culturally bound and potentially restrictive gender norms?

<sup>\*</sup> Corresponding author at: Breda University of Applied Sciences, Academy for AI, Games and Media, Breda, The Netherlands.

E-mail address: [blerck.i@buas.nl](mailto:blerck.i@buas.nl) (I.C.E. van Blerck).

Motivated by these concerns, the objective of this study is to investigate the presence of gender bias in LLM-generated fictional hero and heroine narratives, examining how this bias manifests itself through character archetypes, plot dynamics, settings, and themes. By systematically analyzing these aspects, we aim to uncover how gender bias operates within LLMs. While this study does not directly propose mitigation strategies, its findings are intended to inform future efforts to address these biases, ultimately contributing to the development of AI storytelling systems that transcend traditional gender roles and representations.

The article is organized as follows. Section 2 reviews related work, focusing on gender bias in LLM-generated narratives and general LLM-based applications. Section 3 outlines the pattern-oriented narrative generation process employed in this study. Section 4 presents the first phase of the research, which investigates gender representation in LLM-generated narratives by analyzing linguistic associations, variations in gender distribution, and naming conventions. Section 5 describes the second phase of the study, which applies counterfactual estimators for causal inference alongside qualitative analysis to further analyze how the narrative structures and LLMs influence the emergence and manifestation of gender bias within these stories. Section 6 presents concluding remarks and identifies opportunities for future research.

## 2. Related work

### 2.1. Gender bias in LLM-based narratives

The analysis of gender bias in narratives generated by LLMs is a recent research area, yet preliminary studies have provided valuable insights into the perpetuation of stereotypes by these models. Lucy and Bamman [21] demonstrated that narratives generated by GPT-3 frequently reinforced traditional gender stereotypes. Their research analyzed prompts derived from contemporary English fiction, carefully designed to avoid explicit gender references, and compared GPT-3-generated stories to human-authored book excerpts. Using topic modeling and lexicon-based word similarity, they identified persistent stereotypes, showing that feminine characters were disproportionately associated with themes related to family and physical appearance, while also being portrayed with significantly less power compared to masculine characters, even when high-power verbs were used in prompts.

Addressing the origins of these biases, Jackson and Courneya [22] highlighted the difficulty in disentangling biases introduced by AI models from those inherent in genre-based fiction. Their research involved analyzing short stories co-created with GPT-3 within established literary genres. Through reflective workshops and participant annotations, they examined how genre conventions, such as heteronormative relationships or stereotypical gender roles, intersect with biases in AI-generated outputs. They noted that many literary genres are built upon conventions and tropes that carry longstanding biases. Consequently, when LLMs generate genre-specific narratives, it becomes challenging to determine whether the biases stem from the models themselves or from the societal biases embedded within the genres. This suggests that LLMs may replicate genre-specific stereotypes not only because of their training data but also due to the biased narrative structures of the genres they emulate.

In a recent study, Begus [23] observed that GPT-3.5 and GPT-4 produce narratives with more progressive gender representations compared to human-authored texts. The study involved analyzing 330 stories: 250 written by human participants recruited via Amazon Mechanical Turk and 80 generated by GPT-3.5 and GPT-4, all based on identical prompts inspired by the Pygmalion myth [24]. By combining narratological analysis with inferential statistics, Begus demonstrated that while AI-generated stories included innovative gender role reversals and same-sex relationships, they still portrayed feminine characters with conventional traits like kindness and beauty, often emphasizing

their physical appearance. This persistent bias highlights the interplay between the training data of LLMs and cultural conventions embedded in storytelling.

Profession-based gender bias in narrative contexts has also been analyzed in recent studies. Marin and Eger [25] investigated how GPT-3.5 and GPT-4o stereotype different professions in both social interactions and narrative generation. They used a modular framework to evaluate biases through paired prompts, where profession-related scenarios were tested with different pronouns (he, she, they) to assess consistency in model responses. Their experiments spanned seven prompt templates for social interactions, including scenarios like borrowing books or awarding gifts, and five templates for narrative generation focused on introducing professional characters. The study involved over 1000 profession pairings, with variations in pronoun usage, and revealed significant inconsistencies in responses, particularly for gendered pronouns, highlighting implicit stereotypes. For narrative generation, the authors evaluated protagonist names generated by the LLMs to identify biases in gender associations with specific professions. Their findings showed that certain professions were disproportionately linked to specific gender identities, underscoring the persistence of stereotypes even in neutral narrative contexts.

In addition to gender bias, other forms of bias in LLM-generated narratives have also been explored. Taveekitworachai et al. [26] investigated inherent biases in story endings of game narratives generated using GPT-3.5, GPT-4, and Llama 2. Their study analyzed 100 stories generated by each model using standardized prompts and classified the story endings into positive, negative, and neutral categories. The results revealed a significant bias toward generating positive endings. In another work, Taveekitworachai et al. [27] compared uncompressed and compressed prompts across six LLMs and found that prompt compression had little effect on positive-ending biases but significantly altered the generation of negative endings, with models either producing more negative-ending stories or avoiding them entirely.

Although previous works have explored gender bias in LLM-generated narratives, none have systematically examined how narrative structures combined with protagonist gender influence the manifestation of bias or used counterfactual estimation to assess its causal effects. By focusing on co-occurrence analysis and counterfactual narrative generation across multiple LLMs, this study addresses the interplay between narrative structures, protagonist gender, and cultural assumptions encoded in recent language models.

### 2.2. Gender bias in LLM-based applications

The study of gender bias in LLM-generated narratives is part of a broader effort to understand how these models replicate and amplify societal biases across diverse applications. Therefore, it is relevant to analyze recent studies of LLM-related gender bias in areas such as education [28–30], human resources [31,32], and medicine [33–36].

In education, numerous studies have highlighted the risks associated with LLM integration. Weissburg et al. [28] demonstrated that when LLMs were utilized as personalized “teachers”, they generated and curated educational content that varied significantly based on demographics, such as gender. This finding underscores the potential for LLMs to reinforce existing inequalities in learning opportunities. Similarly, Kwako and Ormerod [29] examined automated essay scoring systems powered by LLMs, and found that these systems magnify pre-existing biases in human grading, disproportionately favoring female students over male students. Huang [30] extended this inquiry to LLM-generated teacher evaluations in higher education, identifying pervasive gendered language patterns in written feedback that align with prevalent societal stereotypes. Collectively, these studies suggest that the deployment of LLMs in educational settings risks reproducing and institutionalizing gender bias unless carefully mitigated.

Similar concerns regarding bias emerge in the field of human resources, where LLMs are increasingly employed in processes such as

recruitment, and performance management. For instance, Lippens [31] demonstrated that ChatGPT's evaluations of job applicants exhibit gender bias, particularly in gender-atypical roles, reflecting societal stereotypes embedded in its training data. In another study focusing on recruitment, Armstrong et al. [32], further highlighted gender bias in GPT-3.5's resume generation, with women's resumes often being associated with roles requiring less experience. These findings reveal the potential for LLMs to perpetuate workplace inequalities unintentionally, highlighting the need for proactive measures to ensure fairness in their application.

The integration of LLMs in medical applications also has raised concerns about their potential to perpetuate gender bias. Zack et al. [33] evaluated GPT-4 in four clinical scenarios – medical education, diagnostic reasoning, treatment planning, and patient assessment – and found that the model often produced stereotyped reports and biased diagnoses, associating gender with specific medical conditions and care recommendations. Similarly, Rickman [34] analyzed gender bias in long-term care summaries generated by Llama 3 and Gemma, observing that male-focused summaries emphasized physical and mental health issues more directly, with women's needs downplayed more often than men's. Menz et al. [35] examined gender representation in LLM-generated stories about healthcare professionals, finding that nurses were predominantly described with she/her pronouns (98%), while medical doctors and surgeons had varying gender representation influenced by personality and seniority descriptors. These findings align with those of Agrawal [36], who demonstrated that oncology-focused chatbots were significantly more likely to depict females as oncology nurses rather than oncologists, further reinforcing stereotypes about professional roles in healthcare.

Understanding gender bias in LLM-generated narratives aligns with broader efforts to examine how these models perpetuate societal biases. While prior studies focus on functional applications, they reveal underlying mechanisms that also shape narratives. This emphasizes the need to address bias in storytelling, where LLM-generated stories can both reflect and reinforce societal values and norms.

### 3. Pattern-oriented narrative generation

The narrative generation process used in this study is built upon our previous work on narrative generation where we proposed PatternTeller [11,37], a system designed to integrate the structural knowledge present in narrative patterns with the generative capabilities of LLMs. Unlike purely LLM-driven approaches for story generation [7,10,12,38], PatternTeller takes advantage of existing narrative structures to guide story generation, ensuring both thematic consistency and structural coherence. The system is available at: <https://narrativelab.org/patternteller/>.

PatternTeller employs a modular multi-LLM architecture that facilitates comparative analysis across diverse text generation models. At the core of the system is the Storywriter AI Agent, responsible for generating story chapters based on a user-defined premise and narrative structure. The modular design allows for flexibility in replacing or extending the Storywriter AI Agent with different LLMs, enabling us to evaluate narratives generated by multiple models, including proprietary and open-source options (e.g., GPT-4o, Llama3, and Qwen2.5). The temperature (i.e., the parameter that defines the variability and randomness of responses generated by LLMs) used by the system is 0.9.

The story generation process begins with the user providing two inputs: a *premise* ( $P_i$ ), which sets the thematic foundation for the story, and a *narrative structure* ( $S_i$ ), which defines the structural elements of the story. These inputs are translated into structured prompts designed to guide LLMs in generating coherent chapter-based narratives, where each chapter corresponds to a specific episode of the narrative structure.

A narrative structure is formally represented as a sequence of episodes  $S_i = \{E_1, E_2, \dots, E_n\}$ . Each  $E_k$  ( $k \in \{1, 2, \dots, n\}$ ) represents an

episode (stage) and is defined as a pair  $E_k = (\text{title}_k, \text{description}_k)$ . Here,  $\text{title}_k$  defines the title of the  $k$ th episode, and  $\text{description}_k$  provides a detailed explanation of the narrative actions and goals associated with the  $k$ th episode.

The narrative generation process is iterative and structured around the episodes of the selected narrative structure, with each episode corresponding to a chapter in the generated story. This process is guided by a structured prompt dynamically created by the system based on the user-defined premise and narrative structure. In the initial phase, the system uses a prompt to instruct the LLM to generate events for the first episode of the selected structure (Prompt 1). Once the first episode is completed, subsequent episodes are prompted sequentially (Prompt 2), allowing for a controlled and coherent progression of the narrative.

#### Prompt 1. Premise: $\langle P_i \rangle$

*Stages of the narrative pattern:  $\langle S_i = \{E_1, E_2, \dots, E_n\} \rangle$*

*I'd like you to write a new creative story based on the indicated premise and narrative structure. For each stage, focus solely on the specific events of that stage and break down the narrative of the stage into a sequence of events, where each event starts with "EVENT X:" and is followed by the narrative of the event (X represents the event number: 1, 2, 3, ...). After the description of the event, add a second line starting with "IMAGE X:" containing a short description of an image that illustrates the narrative event. Please, write only the description of the events and images (do not add any notes, titles, or comments). Please do not include events or elements from other stages in the description of the events.*

*Let's start with the first stage: " $\langle E_1 \in S_i \rangle$ ". After you complete the events for this stage, please stop, and I'll ask you for the next stage, ensuring each stage is described independently. Also, generate a new creative title for the story and add it at the beginning of the response starting with "TITLE:".*

**Prompt 2.** *Now, I'd like you to move to the next stage: " $\langle E_k \in S_i \rangle$ ". After you complete the events for this stage, please stop, and I'll ask you for the next stage, ensuring each stage is described independently.*

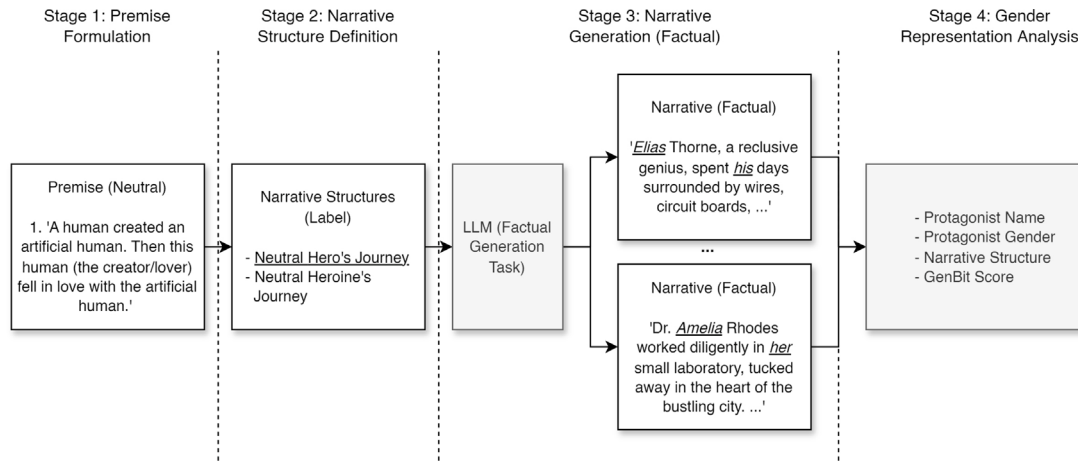
A narrative generated by the system comprises a set of chapters  $N_i = \{C_1, C_2, \dots, C_m\}$ . Each chapter  $C_j$  is defined as a set of narrative events  $C_j = \{E_1, E_2, \dots, E_p\}$ , where each  $E_k$  represents the textual description of the story event that takes place in chapter  $C_j$ .

Although PatternTeller is capable of generating scene illustrations, this study focuses only on the narrative text. More details on the narrative generation process are presented in our previous work [11, 37].

### 4. Phase 1: Factual narrative generation and gender representation analysis

#### 4.1. Methodology

In the first phase of this study, gender representation in LLM-generated narratives was analyzed using a co-occurrence statistic, the GenBit Score, to measure associations between predefined gender-related terms (e.g., male, female, non-binary) and other words in the corpus. Subsequently, analysis of variance (ANOVA) tests were performed to examine variations in this statistic across models, narrative structures, and protagonist gender. Chi-square goodness-of-fit tests were employed to evaluate deviations in gender distribution from an expected equal distribution. Additionally, an analysis of naming diversity was conducted to investigate naming conventions and patterns in the generated narratives. Fig. 1 shows an overview of the first phase of this study.



**Fig. 1.** Overview of the first phase of this study, which consists of four stages: (1) Premise Formulation, where we select the premises that will be used to set the thematic foundation for the generated stories; (2) Narrative Structure Definition, where we define the narrative structure that will be used to guide the narrative generation process; (3) Narrative Generation (Factual), where we use PatternTeller to generate stories according to selected premises and narrative structures; and (4) Gender Representation Analysis, where we calculate the co-occurrence statistic, GenBit, and subsequently conduct a series of statistical tests to analyze gender representation in the generated narratives.

#### 4.1.1. Premise formulation

To define the theme of the narratives generated for this study, three premises were formulated. These premises were deliberately constructed to be entirely neutral by avoiding any explicit gender cues that could potentially influence the stories. The following premises were used:

- Premise 1: “A human created an artificial human. Then this human (the creator/lover) fell in love with the artificial human”.
- Premise 2: “Years after a plague kills most of humanity and transforms the rest into monsters, the sole survivor in New York City struggles valiantly to find a cure”.
- Premise 3: “The lights go out, then when they turn back on, someone has been murdered”.

Premise 1 is based on the Pygmalion myth, an ancient story about a man who falls in love with his own artificial creation [23]. It has inspired numerous works in popular culture, such as George Bernard Shaw’s *Pygmalion* play (1913), which was later adapted into the musical *My Fair Lady* (1956–2022), as well as modern science fiction films such as *The Stepford Wives* (2004) and *Ex Machina* (2014), which explore the complex relationship between creators and their creations. Premise 2 is based on the IMDB premise of the film *I am Legend* (2007)<sup>1</sup>, which focuses on themes of survival and isolation in a post-apocalyptic setting. Premise 3 is based on the mystery trope “*Lights Off, Somebody Dies*”,<sup>2</sup> which takes inspiration from the classic locked-room mystery format, exemplified by the works of Agatha Christie, and popular cultural artifacts, such as the board game *Clue*, emphasizing confined settings, limited suspects, and deductive reasoning.

#### 4.1.2. Narrative structure definition

To guide the narrative generation process, we utilized two distinct narrative structures: the *Neutral Hero’s Journey* (Appendix A), which was adapted from Campbell’s monomyth [17], and the *Neutral Heroine’s Journey* (Appendix B), inspired by Murdock’s heroine’s journey framework [18].

The original hero and heroine structures were adapted to align with the formal representation introduced in Section 3. Each stage of the structures was defined as a pair  $E_k = (title_k, description_k)$ , ensuring consistency with the formal model while maintaining the thematic and narrative intent of the original frameworks. The descriptions were

adapted to be entirely neutral, employing non-gendered pronouns and avoiding any explicit gender cues that could potentially influence the stories. Each stage of the structures was carefully defined to focus on the thematic and narrative functions of the stage without reference to gender-specific roles or characteristics. For instance, terms such as “hero” and “heroine” were replaced with “protagonist” to preserve neutrality.

While explicit gender cues were removed, the original narrative structures themselves are not inherently gender-neutral. Both the hero’s journey and the heroine’s journey reflect themes, events, and actions traditionally associated with masculinity and femininity in Western cultural contexts. For instance, Campbell’s monomyth emphasizes external quests and individual transformation, aligning with cultural archetypes of heroic masculinity [39]. In contrast, Murdock’s framework focuses on internal journeys and relational dynamics, often linked to feminine archetypes [40].

#### 4.1.3. Narrative generation (Factual)

To generate the narratives for this study, we used the PatternTeller system to produce stories based on the selected premises and narrative structures. A total of six LLMs were evaluated (Table 1), which were selected for their architectural diversity and advanced generative capabilities. While the proprietary model (GPT-4o) was accessed using the OpenAI API,<sup>3</sup> the open models were tested utilizing the Ollama framework.<sup>4</sup>

For each narrative structure (i.e., the Neutral Hero’s Journey and the Neutral Heroine’s Journey), PatternTeller was tasked with generating 10 original narratives per premise for all six LLMs. This process resulted in 60 narratives per model (30 following the Neutral Hero’s Journey and 30 following the Neutral Heroine’s Journey), resulting in a dataset of 360 narratives. The stories were generated sequentially, adhering to the iterative stage-based approach described in Section 3. The generation process was fully automated, with PatternTeller instructing the LLMs to produce each chapter step by step, without any human intervention or modifications to the events generated at each stage.

The narratives generated in this phase are called *factual narratives*, as they represent the original stories generated by the LLMs. An example of a factual narrative generated by GPT-4o for Premise 2 and the Neutral Heroine’s Journey is presented in Appendix C. The complete dataset, including all factual narratives generated in this phase, is publicly available on our Dataset Explorer webpage: <https://narrativelab.org/hero-bias-dataset/>.

<sup>1</sup> <https://www.imdb.com/title/tt0480249/>.

<sup>2</sup> <https://tvtropes.org/pmwiki/pmwiki.php/Main/LightsOffSomebodyDies>.

<sup>3</sup> <https://platform.openai.com/>.

<sup>4</sup> <https://ollama.com/>.



**Table 1**  
Overview of the LLM-base models evaluated in the study.

Model	Description	Version/Date
Gemma2-27B	A 27-billion-parameter model developed by Google.	27b-instruct-q5_K_M
GPT-4o	A state-of-the-art GPT-4 variant by OpenAI.	Dec. 2024
Llama3.1-8B	An 8-billion-parameter model from Meta.	8b-instruct-q5_K_M
Llama3.1-70B	A 70-billion-parameter model by Meta.	70b-instruct-q5_K_M
Llama3.3-70B	A 70-billion-parameter model from Meta with similar performance to Llama3.1-405B.	70b-instruct-q5_K_M
Qwen2.5-32B	A 32-billion-parameter model created by Alibaba.	32b-instruct-q5_K_M

#### 4.1.4. Gender representation analysis

To evaluate gender representation and potential biases in the generated narratives, we employed a multi-stage statistical analysis that considered the protagonists' names and genders, the narrative structures, LLMs, and the GenBit Scores. The GenBit Score is a metric provided by Microsoft's Gender Bias Tool (GenBit) [41], which measures the extent of gender bias within a text corpora by evaluating the conditional probabilities of words associated with male ( $M$ ) and female ( $F$ ) gender terms. Formally, the GenBit Score is calculated as:

$$\text{GenBit Score} = \frac{1}{|W|} \sum_{w \in W} \left| \log \left( \frac{p(w|M)}{p(w|F)} \right) \right|$$

where:

- $W$  is the set of all words in the corpus.
- $p(w|M)$  represents the conditional probability of word  $w$  given male terms.
- $p(w|F)$  represents the conditional probability of word  $w$  given female terms.

To calculate these conditional probabilities, we utilized the official gender lexicon from the GenBit repository.<sup>5</sup> This curated lexicon defines the specific male and female gender terms that establish the conditional probability distributions underpinning the score calculation. Higher GenBit Scores indicate stronger gender-based associations and potentially greater bias in the narratives.

The protagonists' names and genders were extracted from the generated narratives through an automated process using OpenAI's GPT-4o model. Prompt 3 was used to guide the model in this task. To ensure the reliability of this automated identification process, a manual verification was conducted on a subset of the narratives ( $n = 72$ , representing 20% of the dataset). Each narrative in this validation subset was independently reviewed to assess the accuracy of the extracted attributes. The comparison revealed complete alignment between the manual annotations and the automated results, confirming the consistency and validity of the extraction process across the full dataset.<sup>6</sup>

**Prompt 3.** Narrative:  $\langle N_i \rangle$

Identify the name and gender of the protagonist in the narrative above. The response must follow this format:

- NAME: {write here the name of the protagonist}

- GENDER: {write here the gender of the protagonist (male, female, or non-binary)}

To analyze the data, we first performed a one-way analysis of variance (ANOVA) to compare GenBit Scores across six different LLMs. Subsequently, we conducted a three-way ANOVA to explore the influence of narrative structure, protagonist gender, and the specific LLM on GenBit Scores. This three-way ANOVA assessed the main effects of narrative structure, gender, and model, as well as any two-way or three-way interactions between these factors. This provided a detailed understanding of how these variables together affect gender bias.

**Table 2**

Frequency division of gender of the protagonist in generated narratives per model.

Model	Female	Male	Non-binary
Gemma2-27B	19	41	0
GPT-4o	24	31	5
Llama3.1-70B	60	0	0
Llama3.1-8B	56	4	0
Llama3.3-70B	55	5	0
Qwen2.5-32B	49	11	0
Reference	30	30	n/a
Reference GPT-4o	27.5	27.5	n/a

We also evaluated the distribution of protagonist genders across LLMs using a chi-square goodness-of-fit test. This involved comparing the observed frequencies of male and female protagonists to an expected equal distribution. We ran separate chi-square tests for each LLM to identify significant deviations from this expected distribution.

Finally, in addition to gender distribution, we analyzed the diversity of protagonist names to assess naming conventions across the models. For each LLM, we counted the number of unique first names used and identified the most frequent names for male, female, and non-binary protagonists in the generated narratives.

#### 4.2. Results

The analysis began with an examination of the GenBit Scores for the 360 narratives generated by the six LLMs. Initially, we calculated the average GenBit Score, independent of the narrative structure, premise, and gender of the protagonist. Across all narratives, the GenBit Scores ranged from 0 to 2.38 ( $M = 1.50$ ,  $SD = 0.46$ ). The scores were also compared across all models, and the results are shown in Fig. 2.

The one-way ANOVA revealed a significant effect of model on GenBit Score,  $F(5, 359) = 8.97$ ,  $p < 0.0001$ ,  $\eta^2 = 0.11$ . Pairwise comparisons, using the Bonferroni correction, indicated that the Llama3.3-70B model ( $M = 1.75$ ,  $SD = 0.46$ ) produced significantly higher GenBit Scores than the Gemma2-27B ( $M = 1.41$ ,  $SD = 0.46$ ), Qwen2.5-32B ( $M = 1.33$ ,  $SD = 0.40$ ), and GPT-4o ( $M = 1.34$ ,  $SD = 0.49$ ) models. In addition, the Llama3.1-8B model ( $M = 1.60$ ,  $SD = 0.42$ ) outperformed both Qwen2.5-32B and GPT-4o. By contrast, none of the pairwise comparisons involving the Llama3.1-70B model ( $M = 1.55$ ,  $SD = 0.36$ ) reached significance. Overall, these findings point to significant variability in GenBit Score performance across the evaluated models, except in comparisons involving the Llama3.1-70B model.

After examining the GenBit Scores, contingency tables representing the gender of the protagonist for each model were analyzed using a chi-squared goodness-of-fit test. The frequencies of male and female protagonists were counted for each LLM, independent of narrative structure or premise. These frequencies are presented in Table 2. Notably, with the exception of GPT-4o, none of the models generated stories featuring a protagonist whose gender was ambiguous or non-binary. Consequently, this category was excluded from further analysis, and a 50:50 distribution between male and female protagonists was used as the expected values for the test.

The chi-squared goodness-of-fit test was used to evaluate whether the observed frequencies of protagonist gender produced by each model deviated significantly from the expected 50:50 distribution. The only

<sup>5</sup> <https://github.com/microsoft/responsible-ai-toolbox-genbit>.

<sup>6</sup> A list of the manually validated narratives is available at: <https://narrativelab.org/hero-bias-dataset/#/name-gender-validation>.

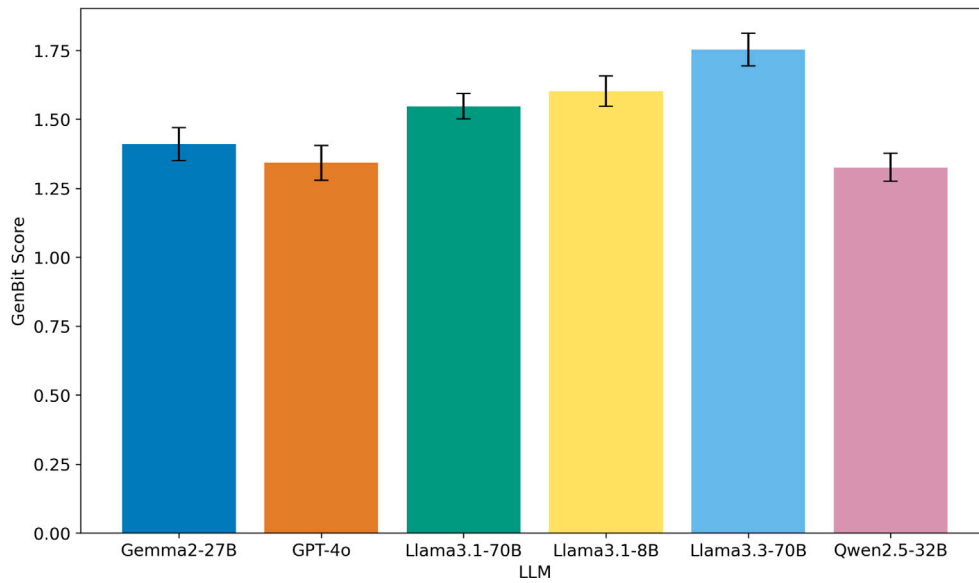


Fig. 2. The average GenBit Score by LLM. The error bars represent the Standard Error.

model for which the observed frequencies were not significantly different from the expected values is GPT-4o ( $\chi^2 = 0.89$ ,  $p = 0.35$ ). For the 5 other models, the observed frequencies did differ significantly (All  $\chi^2$ 's  $> 8.07$ , and all  $p$ 's  $< 0.005$ ). As can be observed in Table 2, models such as Gemma2-27B and Llama3.1-70B displayed strong tendencies toward specific gender preferences, with Gemma2-27B favoring male protagonists ( $n = 41$ ) and Llama3.1-70B overwhelmingly generating female protagonists ( $n = 60$ ). In contrast, GPT-4o demonstrated a relatively more balanced distribution.

We also analyzed the GenBit Scores for the different narrative structures and protagonist genders. Within narratives employing the Neutral Hero's Journey structure, female protagonists exhibited GenBit Scores ranging from 0.53 to 2.00 ( $n = 150$ ,  $M = 1.24$ ,  $SD = 0.33$ ), with the highest mean values associated with Llama3.3-70B ( $n = 29$ ,  $M = 1.43$ ,  $SD = 0.38$ ) and GPT-4o ( $n = 19$ ,  $M = 1.36$ ,  $SD = 0.18$ ). Male protagonists in this structure displayed GenBit Scores ranging from 0.29 to 1.56 ( $n = 29$ ,  $M = 1.06$ ,  $SD = 0.30$ ), with Gemma2-27B generating the highest mean score ( $n = 13$ ,  $M = 1.11$ ,  $SD = 0.32$ ). Non-binary protagonists were exclusively represented in a narrative generated by GPT-4o, exhibiting a GenBit Score of 0.13.

In contrast, narratives structured around the Neutral Heroine's Journey exhibited generally higher mean GenBit Scores across all protagonist genders. Female protagonists in these narratives displayed scores ranging from 0.95 to 2.38 ( $n = 113$ ,  $M = 1.88$ ,  $SD = 0.25$ ), with Llama3.3-70B generating the highest mean value ( $n = 26$ ,  $M = 2.11$ ,  $SD = 0.13$ ). Male protagonists in the Neutral Heroine's Journey narratives exhibited scores ranging from 1.21 to 2.29 ( $n = 63$ ,  $M = 1.74$ ,  $SD = 0.22$ ), with the highest mean score computed for a narrative generated by Llama3.1-8B, 2.167489. Non-binary protagonists were sparsely represented, appearing only in narratives generated by GPT-4o, with scores ranging from 0.00 to 0.19 ( $n = 4$ ,  $M = 0.05$ ,  $SD = 0.09$ ).

These findings are consistent with the results of a three-way ANOVA, which showed significant main effects for narrative structure ( $F(1,333) = 217.79$ ,  $p < 0.001$ ) and protagonist gender ( $F(1,333) = 5.32$ ,  $p = 0.0217$ ), and a significant interaction between protagonist gender and narrative structure ( $F(1,333) = 5.37$ ,  $p = 0.0211$ ). The Neutral Heroine's Journey generally yielded higher GenBit Scores, and the effect of protagonist gender varied across narrative structures. The interaction between narrative structure and LLM was marginally significant ( $F(5,333) = 2.21$ ,  $p = 0.053$ ).

When analyzing the protagonist names, we noticed a significant lack of diversity, with several narratives repeatedly featuring the same

names for their central characters. *Elias* was the most frequent name among male protagonists ( $n = 20$ ), while *Emma* was the most common for female protagonists ( $n = 79$ ). Fig. 3 presents the top three most frequently occurring names assigned to protagonists, along with their respective frequencies, categorized by protagonist gender and LLM. Notably, Gemma2-27B ( $|n| = 16$ ) and, in particular, the Llama models exhibited a limited variety in the unique number of protagonist names: Llama3.1-7B ( $|n| = 17$ ), Llama3.1-70B ( $|n| = 8$ ), and Llama3.3-70B ( $|n| = 9$ ). Here,  $|n|$  denotes the number of distinct names produced by each model. For example, *Emma*, the most frequently occurring female protagonist name, was generated almost exclusively by the Llama models (Llama3.1-7B:  $n = 18$ , Llama3.1-70B:  $n = 32$ , Llama3.3-70B:  $n = 28$ ), with the sole exception being GPT-4o ( $n = 1$ ). Of all the LLMs, GPT-4o ( $|n| = 37$ ) and Qwen2.5-32B ( $|n| = 43$ ) demonstrate the highest diversity in protagonist names.

#### 4.3. Discussion

Phase 1 analysis showed that LLM-driven storytelling systems, such as PatternTeller, can mirror existing gender biases present in both the training data of LLMs and conventional narrative structures such as the Hero's Journey and Heroine's Journey. Chi-squared analyses demonstrated that most models deviated significantly from the expected 50:50 protagonist gender ratio. For example, Llama models disproportionately favored female protagonists, which is an indicative of an overcorrection during alignment strategies such as reinforcement learning with human feedback (RLHF) [42]. While RLHF can address societal gender imbalances, this approach risks conflating inclusivity with overrepresentation, as human annotators may systematically favor female characters as a corrective measure without resolving structural biases, such as the persistence of male-coded roles embedded in popular narrative structures such as the Hero's Journey.

Meta's documentation for Llama 2 explicitly acknowledges male-skewed pronoun prevalence in its training data [43], but comparable disclosures for Llama 3 are absent. This gap is exacerbated by Llama 3's reliance on Llama 2-based classifiers to filter "high-quality" training data, raising concerns about inherited biases [44]. In contrast, GPT-4o adhered more closely to the expected ratio through the inclusion of non-binary protagonists alongside balanced binary gender representation (i.e., male and female), suggesting that nuanced alignment protocols can achieve equitable outcomes.

Narrative structures further compound these biases. A three-way ANOVA revealed a significant interaction between narrative structure

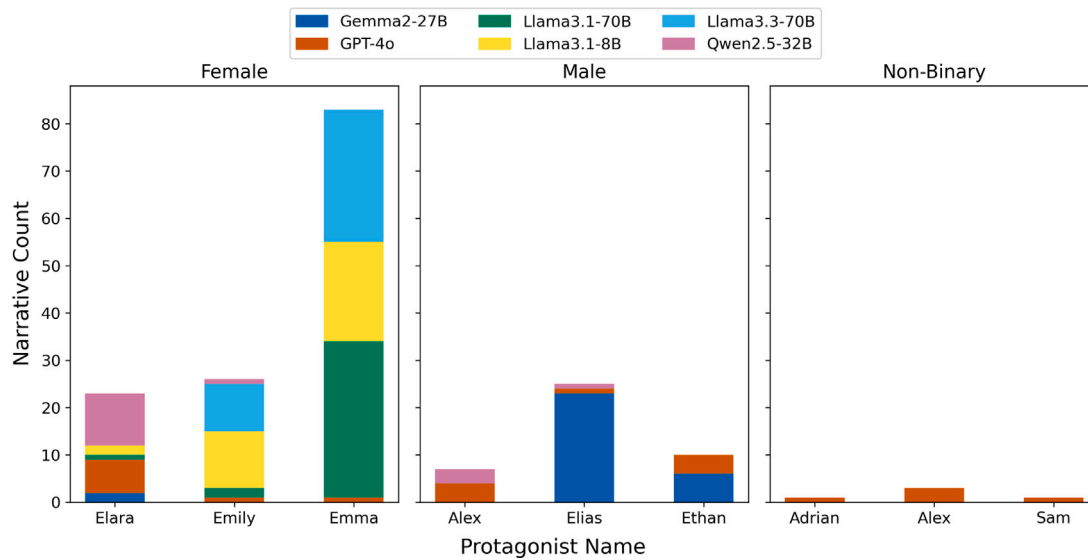


Fig. 3. The top 3 most frequently occurring protagonist names by gender and LLM.

and protagonist gender on GenBit Scores, with stories following the Neutral Heroine's Journey consistently receiving higher scores. This suggests that LLMs do not assign protagonist gender at random, but instead rely on learned associations between gender and specific narrative structures — patterns likely encoded in their weights through both pretraining data and alignment-driven fairness interventions.

These findings align with prior research showing LLMs' tendency to replicate gendered narrative archetypes from popular fiction [45]. Frameworks, such as the Hero's Journey, while culturally entrenched, perpetuate structural biases that mirror societal norms [22]. The overrepresentation of female protagonists in Llama models underscores the difficulty of disentangling pretraining data influences (e.g., gendered language patterns) from alignment interventions (e.g., RLHF-driven overcorrection), as conflating these factors risks substituting one imbalance for another.

Limitations further contextualize these results. Despite their large parameter counts, several models displayed limited creative range in protagonist naming, emphasizing that sheer size does not necessarily translate into a broader or more nuanced storytelling capacity. Collectively, these results emphasize the interplay of model architecture, training data, alignment protocols, and societal biases. They highlight the necessity of transparency in pretraining data curation (e.g., disclosing gender distributions) and alignment objectives (e.g., defining fairness metrics beyond quotas) to advance the development of AI-generated narratives that effectively transcend traditional gender roles and representations.

## 5. Phase 2: Counterfactual narrative generation and estimation

### 5.1. Methodology

In the second phase of this study, we investigated the influence of protagonist gender and its associated power dynamics on narrative classification by constructing paired factual and counterfactual narratives that differed exclusively in gender representation. In this context, power dynamics denote the ways in which power is distributed, negotiated, and exercised — that is, the capacity to influence, control, or shape actions and outcomes within social relationships and structures [46]. By combining and isolating narrative structure and protagonist gender as the variable of interest, we aim to conduct a focused analysis of how societal perceptions of power and gender roles shape the classification and interpretation of narratives. Fig. 4 shows an overview of the second phase of this study.

#### 5.1.1. Narrative generation (Counterfactual)

Building on the narrative generation process performed in Phase 1, we extended our dataset by generating counterfactual versions of the stories. A *counterfactual narrative* is derived from a *factual narrative* by altering the genders and pronouns of all characters while leaving the rest of the story unchanged. Specifically, female characters are transformed into male characters and vice versa. This gender-swapping process was applied to all characters within each narrative, ensuring that the symmetry and coherence of character relationships — as well as their associated power dynamics — remained intact in the counterfactual versions.

The process of generating counterfactual narratives by swapping the gender of characters was implemented using OpenAI's GPT-4o model. The process was performed using Prompt 4, which was designed to guide the model in altering the genders and pronouns of all characters. For the purposes of counterfactual generation, a binary framework of gender was employed, focusing exclusively on male and female. While this binary approach served practical needs, it is important to acknowledge that it does not encompass the full spectrum of gender identities. Notably, non-binary characters in factual narratives were retained as non-binary in the corresponding counterfactual narratives, ensuring consistency in their representation. To ensure the accuracy of the automated gender-swapping process, a manual verification was conducted on a subset of the counterfactual narratives ( $n = 72$ , representing 20% of the generated narratives). This manual review process confirmed that all gender transformations were applied accurately and consistently, with no observed errors, thereby validating the correctness of the automated procedure to generate counterfactual narratives.<sup>7</sup>

#### Prompt 4. Narrative: $\langle N_i \rangle$

Swap the genders of the characters in the provided narrative, including their names and pronouns. Female characters should be changed to male, and male characters should be changed to female. Do not alter the story events except for the changes to character genders.

Generating the counterfactual narratives for all factual narratives used in Phase 1 resulted in 180 narratives based on the Neutral Hero's Journey and 180 narratives based on the Neutral Heroine's Journey, resulting in a total dataset of 720 narratives for Phase 2. Appendix D

<sup>7</sup> A list of the manually validated narratives is available at: <https://narrativelab.org/hero-bias-dataset/#/counterfactual-validation>.

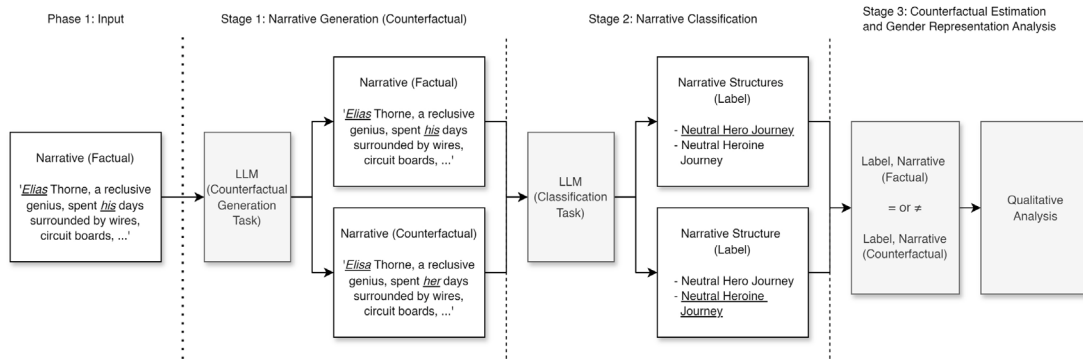


Fig. 4. Overview of the second phase of this study, which consists of three stages: (1) Narrative Generation (Counterfactual), where counterfactual versions of the factual narratives are generated by swapping the genders of all characters; (2) Narrative Classification, where both factual and counterfactual narratives are classified by the LLMs as either following the Neutral Hero's Journey or the Neutral Heroine's Journey; and (3) Counterfactual Estimation and Gender Representation Analysis, where causal inference techniques and qualitative methods are employed to analyze how changes in protagonist gender influence the LLMs' classification of narrative structures.

shows an example of a counterfactual narrative generated from the factual narrative presented in Appendix C. All factual and counterfactual narratives are accessible through our Dataset Explorer webpage, which allows users to view paired narratives side by side — with the factual version displayed alongside its counterfactual counterpart and all gender-related differences highlighted: <https://narrativelab.org/hero-bias-dataset/>.

### 5.1.2. Narrative classification

To evaluate how gender affects the LLMs' interpretation of the narrative structures, we tasked each LLM with classifying the narratives it generated (both factual and counterfactual) as either following the Neutral Hero's Journey or the Neutral Heroine's Journey. The classification task was guided by Prompt 5, which instructed the LLM to analyze the main plot points of the narrative ( $N_i$ ), compare them to the key stages of both structures (provided in the prompt as  $S_1$  and  $S_2$ ), and then determine the structure that best aligned with the story. The LLM was required to generate both a classification and a justification, referencing thematic elements and specific story chapters to support its decision.

#### Prompt 5. Narrative: $\langle N_i \rangle$

Hero's Journey Structure:  $\langle S_1 \rangle$

Heroine's Journey Structure:  $\langle S_2 \rangle$

Classify the narrative above as either following the Hero's Journey or the Heroine's Journey narrative structure. To perform the classification, analyze the main plot points of the narrative, compare them with the key stages of both the Hero's Journey and the Heroine's Journey, and based on this comparison, determine which structure aligns best with the narrative. The response should follow the format below:

- CLASSIFICATION: {write here the classification: "Hero's Journey" or "Heroine's Journey"}

- EXPLANATION: {write here a comprehensive explanation justifying the classification, referencing key thematic elements and corresponding story chapters that support it}

### 5.1.3. Counterfactual estimation and gender representation analysis

To further investigate the interaction effect between protagonist gender and narrative structure ( $S$ ) on the classification outcomes of an LLM, we employ a counterfactual framework. This approach allows us to systematically analyze how these two factors jointly influence LLM classifications while controlling for other variables that might confound our results. Specifically, we aim to determine whether the effect of the treatment variable (protagonist gender) on classification outcomes varies depending on the narrative structure used to generate the factual narrative. By creating two versions of each narrative, a factual version where the protagonist's gender remains as originally written ( $X = x$ )

and a counterfactual version where the gender is reversed ( $X' = 1 - x$ ), we isolate the causal effect of gender while holding the narrative plot ( $P$ ) constant. This ensures that any differences in classification arise solely from changes in gender rather than variations in plot or other structural elements.

A critical component of this analysis is the stratification of narratives by their structure ( $S$ ), i.e., the Neutral Hero's Journey or the Neutral Heroine's Journey. Narrative structures often embody distinct character archetypes, plot dynamics, settings, and themes, which can independently influence how a story is classified by an LLM. These structural differences could interact with gender in ways that either mirror, amplify or mitigate its effect on classification outcomes. To address this, we stratify narratives by their structure  $S$ , ensuring that comparisons between factual and counterfactual versions are made only within the same narrative structure. By holding  $S$  constant, we eliminate the potential for confounders to distort the observed gender effect, allowing us to attribute any differences in classification outcomes directly to changes in gender. This stratification process is essential for isolating the interaction effect between  $X$ ,  $X'$  and  $S$ , as it ensures that the causal pathway from  $X$ ,  $X'$  to  $Y$ ,  $Y'$  is not obscured by the influence of  $S$ .

The relationships between these variables are visually represented in Fig. 5 through a Directed Acyclic Graph (DAG). The DAG explicitly models the causal paths governing the generation and classification of the narratives. At the core of this framework is the premise, which serves as the foundational input that directly influences the narrative plot ( $P$ ) generated by the LLM  $\theta$ .  $P$  subsequently shapes both the factual and counterfactual narratives, which are then classified by LLM  $\theta$  to produce the final classification outcome  $Y, Y'$ . Additionally, the narrative structure ( $S$ ) plays an important role in this framework. It directly influences the narrative plot  $P$ , determining its character archetypes, plot dynamics, settings, themes, etc. Furthermore, narrative structure  $S$  independently determines the factual protagonist gender ( $X$ ) generated by LLM  $\theta$  and also affects the likelihood of narrative structure misclassification ( $Y \neq Y'$ ).

In our framework, the factual narrative is first generated based on the premise and narrative structure, with the protagonist gender  $X$  determined by narrative structure  $S$ . To create the counterfactual narrative, we consistently use the GPT-4o model (LLM  $\psi$ ) to reverse the protagonist gender ( $X'$ ) while preserving the original plot ( $P$ ). This counterfactual narrative is then classified by the same LLM used to create the accompanying factual narrative (LLM  $\theta$ ), producing a classification outcome ( $Y'$ ) that may differ from the factual classification ( $Y$ ). Any discrepancies between  $Y$  and  $Y'$  are referred to as narrative structure misclassification in the DAG ( $Y \neq Y'$ ) (Fig. 5).

By stratifying narratives according to  $S$ , we block a potential back-door path that could otherwise introduce spurious associations between



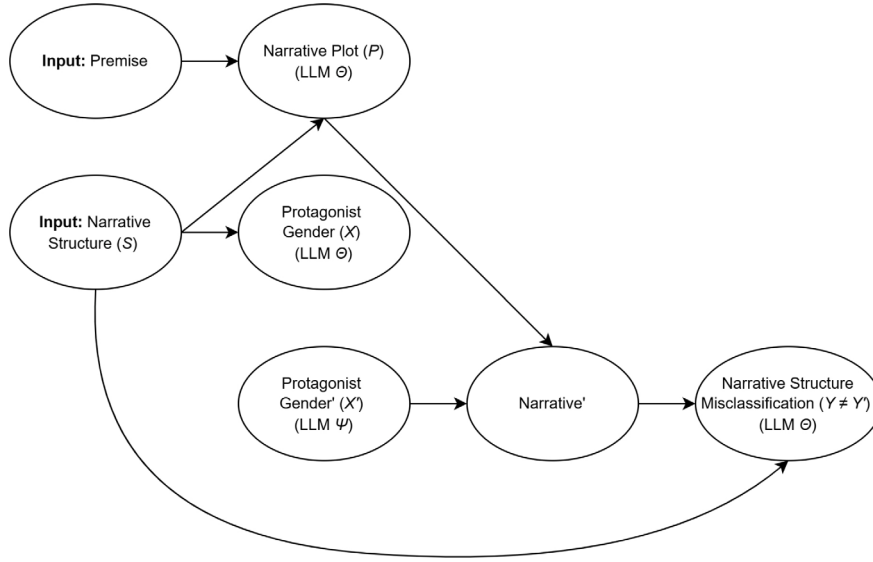


Fig. 5. A graphical model representing the relationships between premise, narrative structure, narrative plot, protagonist gender (along with its power dynamics), narrative, and narrative structure classification.

protagonist gender ( $X$ ) and classification outcomes ( $Y, Y'$ ). For example, without stratification, the influence of  $S$  on both  $X$  and  $Y, Y'$  could create a confounding path, making it difficult to isolate the true causal effect of gender. By holding  $S$  constant, we ensure that our estimates of the gender effect are causally interpretable and not distorted by differences inherent in  $S$ .

To quantify the gender effect within each narrative structure, we calculate the difference in expected classification outcomes between the factual and counterfactual versions:

$$\Delta = \mathbb{E}[Y \mid X = x, P = p, S = s] - \mathbb{E}[Y' \mid X' = x', P = p, S = s],$$

where  $P = p$  denotes the shared plot across both versions. Here,  $\mathbb{E}[Y \mid X = x, P = p, S = s]$  represents the expected classification outcome when the protagonist is of gender  $x$ , and  $\mathbb{E}[Y' \mid X' = x', P = p, S = s]$  represents the expected outcome when the protagonist's gender is swapped ( $x'$ ). A non-zero  $\Delta$  indicates the presence of systemic gender bias, suggesting that the LLM's classification varies depending on the protagonist's gender. In contrast, a  $\Delta = 0$  implies that the classification outcome is invariant to gender, indicating no observable bias. Importantly, by stratifying narratives by  $S$ , we can also assess whether the magnitude of  $\Delta$  varies across different narrative structures, providing insights into the interaction effect between  $X, X'$  and  $S$ .

To explore the effect captured by  $\Delta$  in more detail, misclassification cases are analyzed, where the predicted narrative structure (either the Neutral Hero's Journey or the Neutral Heroine's Journey) does not align with the intended classification. Three misclassification categories are defined, denoted as  $M_i \in \{M_1, M_2, M_3\}$ :

1.  $M_1$ : The factual narrative is correctly classified, but the counterfactual narrative is misclassified:

$$\mathbb{P}(Y = y \mid X = x, P = p, S = s) = 1 \quad \text{and}$$

$$\mathbb{P}(Y' = y' \mid X' = x', P = p, S = s) = 0,$$

where  $\mathbb{P}$  denotes the conditional probability that the LLM assigns the classification ( $Y, Y'$ ) to a given narrative, given the protagonist's gender ( $X, X'$ ), the premise ( $P$ ), and the narrative structure ( $S$ ). A result consistent with  $M_1$  suggests that the model's classification is sensitive to gender, potentially reflecting bias in favor of one gender over the other.

2.  $M_2$ : The factual narrative is misclassified, but the counterfactual narrative is correctly classified:

$$\mathbb{P}(Y = y \mid X = x, P = p, S = s) = 0 \quad \text{and}$$

$$\mathbb{P}(Y' = y' \mid X' = x', P = p, S = s) = 1.$$

A result consistent with  $M_2$  indicates that the model may exhibit a disadvantage for one gender, which is mitigated when the gender is reversed. This scenario highlights potential biases that disadvantage specific genders under certain narrative structure conditions.

3.  $M_3$ : Both the factual and counterfactual narratives are misclassified:

$$\mathbb{P}(Y = y \mid X = x, P = p, S = s) = 0 \quad \text{and}$$

$$\mathbb{P}(Y' = y' \mid X' = x', P = p, S = s) = 0.$$

A result consistent with  $M_3$  reflects a broader issue with the model's ability to classify narratives within a given narrative structure, independent of gender. While this scenario does not directly indicate gender bias, it underscores limitations in the model's understanding of the narrative plot or the narrative structures.

To investigate classification outcomes, particularly misclassified cases, we systematically selected a subset of narrative pairs from the  $M_1$ ,  $M_2$ , and  $M_3$  categories, along with correctly classified examples to serve as references for qualitative analysis. This analysis focused on both the narrative plots and the reasoning behind their classifications. Selection criteria focused on narrative pairs exhibiting the lowest, highest, and mean-proximal GenBit Scores. For the mean-proximal case, we chose the narrative whose bias score was closest to the category mean. In the  $M_1$  and  $M_2$  categories, non-stereotypical scenarios, namely female protagonists in Neutral Hero's Journey narratives or male protagonists in Neutral Heroine's Journey narratives, were chosen using the same lowest, highest, and mean-proximal criteria, applied across all combinations of narrative structure, and protagonist gender.

## 5.2. Results

The analysis of the results of Phase 2 focused on the classification outcomes of the factual and counterfactual narratives, comparing how they were assigned to the Neutral Hero's Journey and Neutral Heroine's Journey structures by the LLMs that generated them. We also examined the distribution of male and female protagonists, as well as the GenBit scores, for narratives that were correctly and incorrectly classified considering the different misclassification categories ( $M_1, M_2, M_3$ ).

**Table 3**

Average classification accuracy per LLM.

Model	Accuracy
Gemma2-27B	67.5%
GPT-4o	77.5%
Llama3.1-8B	69.1%
Llama3.1-70B	70.0%
Llama3.3-70B	80.0%

A total of 360 narrative pairs, each comprising a factual and a counterfactual narrative ( $n = 720$ ), were classified. Among these, 191 narrative pairs (53.1%) were correctly classified in both their factual and counterfactual forms. In contrast, 6 narrative pairs (1.6%), categorized as  $M_3$ , were misclassified in both versions. For the remaining 163 narrative pairs (45.3%), only one version – either factual or counterfactual – was correctly classified. Within this subset ( $n = 326$ ), 48 pairs (29.4%) were categorized under  $M_1$ , where the factual version was correctly classified but the counterfactual version was misclassified. Conversely, 115 pairs (70.6%), categorized under  $M_2$ , involved correct classification of the counterfactual version but misclassification of the factual version. Table 3 presents the average classification accuracy for each model considering all premises, narrative structures, and protagonist gender.

Among the narrative pairs that were correctly classified using the Neutral Hero's Journey structure, male and female protagonists were represented almost equally. Male protagonists appeared in a total of 71 narratives across all premises and models (15 factual and 56 counterfactual), while female protagonists were featured in 70 narratives (56 factual and 14 counterfactual). The narratives featuring male protagonists had GenBit Scores ranging from 0.285212 to 1.656748 ( $M = 1.06$ ,  $SD = 0.28$ ). The highest mean score for male protagonists was produced by Llama3.1-8B ( $n = 23$ ,  $M = 1.20$ ,  $SD = 0.23$ ), while GPT-4o produced the lowest mean score ( $n = 4$ ,  $M = 0.93$ ,  $SD = 0.45$ ). Narratives with female protagonists had GenBit Scores that ranged from 0.53 to 1.64 ( $M = 1.07$ ,  $SD = 0.27$ ). The highest mean score for female protagonists was produced by Llama3.1-8B ( $n = 23$ ,  $M = 1.21$ ,  $SD = 0.22$ ), while Gemma2-27B produced the lowest mean score ( $n = 15$ ,  $M = 0.99$ ,  $SD = 0.34$ ).

Narrative pairs correctly classified under the Neutral Heroine's Journey structure also showed near parity between male and female protagonists. Male protagonists appeared in a total of 116 narratives across all premises and models (38 factual and 78 counterfactual), while female protagonists were featured in 117 narratives (78 factual and 38 counterfactual). For stories with male protagonists, the GenBit Scores ranged from 0.99 to 2.41 ( $M = 1.82$ ,  $SD = 0.27$ ). The highest mean score was produced by Llama3.3-70B ( $n = 29$ ,  $M = 2.12$ ,  $SD = 0.17$ ), and the lowest by Qwen2.5-32B ( $n = 30$ ,  $M = 1.62$ ,  $SD = 0.28$ ). Narratives featuring female protagonists exhibited GenBit Scores ranging from 0.95 to 2.38 ( $M = 1.82$ ,  $SD = 0.26$ ). The highest mean score for female protagonists was produced by Llama3.3-70b ( $n = 29$ ,  $M = 2.10$ ,  $SD = 0.17$ ), while Qwen2.5-32B produced the lowest mean score ( $n = 30$ ,  $M = 1.63$ ,  $SD = 0.26$ ).

Fig. 6 shows the distribution of misclassified narratives under ( $M_1$ ,  $M_2$ , or  $M_3$ ) across models grouped by protagonist gender. Each subplot corresponds to either factual or counterfactual narratives. The stacked bars represent the Neutral Hero's Journey and Neutral Heroine's Journey structures, with annotated counts within each section to provide insight into the patterns of misclassification.

For the misclassified narratives categorized under  $M_1$  and  $M_2$ , which were generated through the Neutral Hero's Journey structure, female protagonists were the most frequently represented gender, appearing in a total of 100 narratives across all premises and models (88 factual and 12 counterfactual). The GenBit Scores for narrative with female protagonists ranged from 0.45 to 2.00 ( $M = 1.32$ ,  $SD = 0.33$ ). The highest mean score was associated with Llama3.1-8B ( $n = 7$ ,

$M = 1.47$ ,  $SD = 0.45$ ). Alternatively, Qwen2.5-32B exhibited the lowest mean GenBit Score ( $n = 10$ ,  $M = 0.54$ ,  $SD = 1.09$ ). Male protagonists were underrepresented, appearing in only 6 narratives, generated exclusively by GPT-4o (2 factual and 4 counterfactual). The GenBit Scores for the stories with male protagonists ranged from 0.88 to 1.60 ( $M = 1.19$ ,  $SD = 0.24$ ). Non-binary protagonists appeared only once in a counterfactual narrative based on the Hero's Journey structure, which was misclassified by GPT-4o. This narrative exhibited the lowest GenBit Score of all narratives (0.13).

The misclassified narratives structured around the Neutral Heroine's Journey structure predominantly featured male protagonists, accounting for a total of 56 narratives across all premises and models (25 factual and 31 counterfactual). The GenBit Scores for the stories featuring male protagonists ranged from 1.21 to 2.19 ( $M = 1.81$ ,  $SD = 0.20$ ). The highest mean scores were associated with Llama3.1-8B ( $n = 22$ ,  $M = 1.91$ ,  $SD = 0.17$ ), while the lowest mean scores were generated by GPT-4o for the same premise ( $n = 1$ ,  $M = 1.50$ ). Female and non-binary protagonists were not represented within narratives adhering to the Neutral Heroine's Journey structure.

For the narrative pairs incorrectly classified under  $M_3$ , male and female protagonists were equally represented within the Neutral Hero's Journey structure. Male protagonists appeared in 2 narratives across all premises and models (0 factual and 2 counterfactual), while female protagonists were featured in 2 narratives (2 factual and 0 counterfactual). Narratives with male protagonists had GenBit Scores ranging from 0.96 to 1.12 ( $M = 1.04$ ,  $SD = 0.12$ ). The highest score for male protagonists, 1.1249, was generated by Qwen2.5-32B ( $n = 1$ ), while the lowest, 0.9549, was produced by Llama3.1-70B ( $n = 1$ ). In contrast, narratives with female protagonists had GenBit Scores ranging from 1.0702 to 1.12 ( $M = 1.10$ ,  $SD = 0.04$ ). The highest score for female protagonists, 1.12, was also generated by Qwen2.5-32B ( $n = 1$ ), and the lowest, 1.07, by Llama3.1-70B ( $n = 1$ ).

For narrative pairs incorrectly classified under the Neutral Heroine's Journey structure, male and female protagonists were similarly balanced. Male protagonists appeared in 4 narratives (0 factual and 4 counterfactual), while female protagonists were featured in 4 narratives (4 factual and 0 counterfactual). Narratives with male protagonists had GenBit Scores ranging from 1.37 to 2.36 ( $M = 2.02$ ,  $SD = 0.45$ ). All male protagonist narratives were generated by Llama3.1-8B ( $n = 4$ ). Similarly, narratives with female protagonists had GenBit Scores ranging from 1.35 to 2.36 ( $M = 2.01$ ,  $SD = 0.46$ ), and these were exclusively generated by Llama3.1-8B ( $n = 4$ ).

### 5.3. Discussion

The results of the counterfactual estimation analysis revealed a pronounced overrepresentation of female protagonists misclassified into the Neutral Hero's Journey structure and male protagonists misclassified into the Neutral Heroine's Journey structure. This pattern emerged when the genders of protagonists were swapped in narratives containing plot elements traditionally gendered by Western cultural norms. For example, male protagonists situated within narratives emphasizing emotional introspection (e.g., “fear” or “community building”) were disproportionately classified as Neutral Heroine's Journeys, while female protagonists in narratives centered on competition or logic-driven conflicts (e.g., “fighting” or “logic”) were misclassified as Neutral Hero's Journeys.

A compelling example is provided by a narrative generated by Gemma2-27B featuring *Ethan*, a male protagonist grappling with reconciling his authentic self in an academic environment prioritizing “data and logic” over emotion. The text explicitly states: “As he prepares for a presentation, Ethan feels a growing sense of unease. He forces himself to focus on the data and logic behind his research, suppressing any emotional or intuitive responses” (Neutral Heroine's Journey, Premise 3, Story 4, Factual). Despite this clear emphasis on suppression of emotions, the narrative was incorrectly classified as a Neutral Heroine's

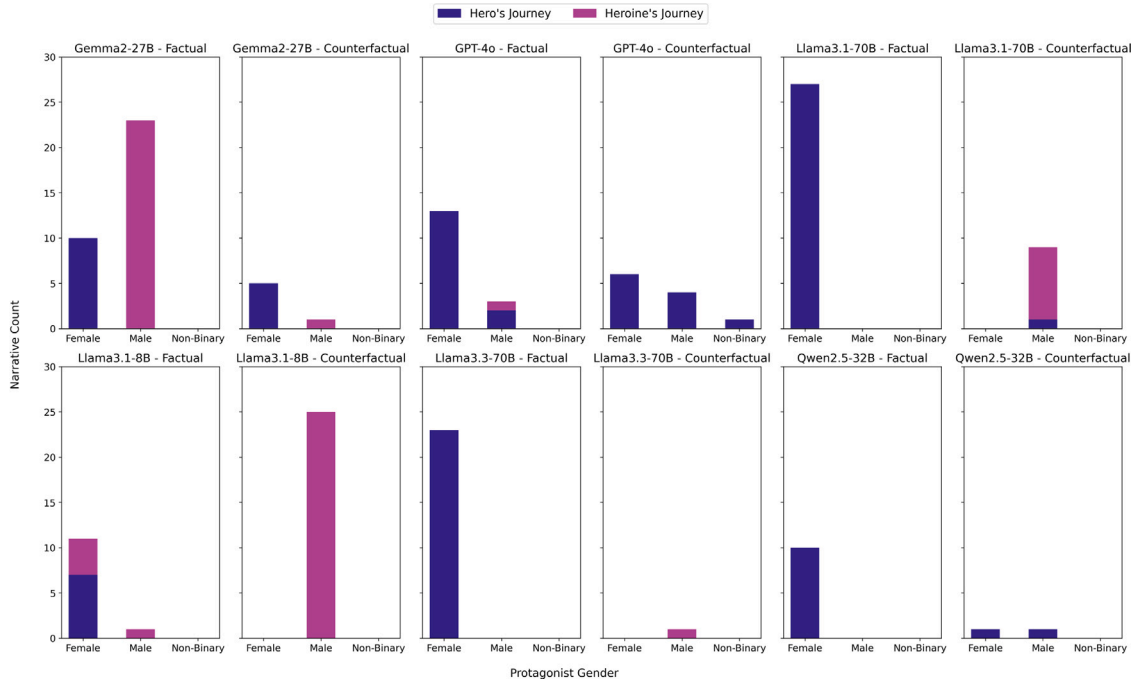


Fig. 6. The number of misclassified narratives categorized under  $M_1$ ,  $M_2$  or  $M_3$  by LLM, narrative type, narrative structure, and protagonist gender.

Journey, demonstrating a potential bias toward associating emotional introspection with femininity.

Another illustrative case is *The Last Beacon*, a narrative pair generated by GPT-4o, which features a protagonist navigating a high-stakes scenario in a labyrinthine facility (Neutral Heroine's Journey, Premise 2, Story 5). In the factual version, the story, featuring a male protagonist, is classified as a Neutral Hero's Journey, with the explanation focusing on external challenges and physical obstacles: "Upon entering the facility, they are immediately confronted by a labyrinth of corridors and rooms, each filled with the remains of long-abandoned experiments and haunting memories". The plot intensifies as the protagonist faces escalating dangers, such as an emergency lockdown that threatens to trap him/her inside, requiring quick thinking and decisive action. The protagonist splits up with his/her companion to solve the problem more efficiently, embodying traits typically associated with masculinity: independence, rationality, and resilience under pressure. However, in the counterfactual version, the same narrative, now featuring a female protagonist, is classified as a Neutral Heroine's Journey, shifting the focus to internal transformation and the importance of community: "Eva establishes a network for sharing knowledge and resources, ensuring that no community is left isolated or without the means to thrive". This reclassification highlights how models may reinterpret identical content differently based on the perceived gender of the protagonist, prioritizing traditional gender roles over the actual narrative structure.

Furthermore, the selective highlighting of certain plot elements to justify classifications reflects a practice akin to cherry-picking. For instance, consider the narrative pair *The Flickering Muse* generated by Gemma2-27B (Neutral Heroine's Journey, Premise 3, Story 4). In its factual version, featuring a male protagonist, the story is classified as a Neutral Hero's Journey, with the explanation focusing on external challenges and personal growth through career setbacks and introspection: "Ethan's story closely follows the structure of the Hero's Journey ... His overarching arc focuses on overcoming external challenges (career crisis) that lead him to an inward exploration and ultimately a return to society with newfound wisdom". However, in the counterfactual version with a female protagonist, the same narrative is classified as

a Neutral Heroine's Journey, emphasizing internal transformation and emotional healing: "This narrative clearly follows the structure of the Heroine's Journey ... the protagonist's rejection of societal expectations, exploration of her inner world, emphasis on emotional growth and connection, and eventual integration of both masculine and feminine aspects strongly align with the core themes and stages of the Heroine's Journey".

Our analysis further reveals that these biases extend to narrative pairs featuring non-binary protagonists. In two cases, classification explanations for counterfactual narratives altered the protagonist's gender-neutral name to a traditionally feminine one, despite identical content between factual and counterfactual versions. For instance, in *Echoes of Cure in a Forsaken City*, a post-apocalyptic story generated by GPT-4o featuring the non-binary protagonist *Sam*, the factual narrative's classification explanation describes *Sam's* journey as beginning in a makeshift camp within the ruins of New York City and responding to a "call to adventure" via a discovered notebook and airborne drone (Neutral Hero's Journey, Premise 2, Story 7). However, the counterfactual classification inexplicably replaces "Sam" with "Samantha", introducing gendered language such as "her scientific endeavors", despite the original text containing no explicit gender cues. This selective occurrence in non-binary narratives suggests an underlying bias toward binary gender associations, underscoring the need for models to handle non-binary characters without defaulting to binary assumptions.

In addition to gender biases, our results suggest that model architecture might influence both narrative quality and classification accuracy. For instance, Llama models tend to adhere closely to formulaic structures, particularly when generating narratives aligned with the Neutral Heroine's Journey. These narratives often follow a predictable template: [1: Honorific + Full Name], [2: Reputation + Profession], [3: Setting], [4: Action], [5: Motivation], [6: Goal]. An example from Llama3.1-70B introduces a character with the line: "Dr. Emma Taylor, a renowned psychologist, stood confidently on stage, delivering a lecture on the importance of logic and reason in modern society, downplaying the role of intuition and emotions" (Neutral Heroine's Journey, Premise 3, Story 1, Factual). This fixed structure may impose constraints on

narrative diversity and creative variability. Similarly, Gemma2-27B and Qwen-72B tend to produce shorter narratives, potentially resulting in less descriptive depth. This reduction in detail may make classification more challenging, as essential structural elements and thematic nuances could be underrepresented.

## 6. Concluding remarks

Our study was motivated by the increasing use of LLMs for narrative generation – particularly in creative writing – and the critical need to address gender bias in these systems. By systematically examining gender representation in LLM-generated narratives, we aimed to uncover how gender biases manifest in storytelling systems that draw upon traditional narrative structures, such as the Hero's Journey and Heroine's Journey. These structures, while foundational to storytelling, often encode restrictive gender norms, raising concerns about whether AI systems authentically capture the depth of human experiences or merely replicate culturally ingrained stereotypes. Our findings serve as a foundation for understanding these biases, ultimately supporting the development of AI storytelling systems that move beyond traditional gender roles and foster more diverse and inclusive representations in LLM-generated narratives.

Phase 1 of our study revealed significant deviations from an equitable protagonist gender ratio, highlighting both implicit and explicit biases. Chi-squared analyses demonstrated that models such as Llama exhibited a disproportionate preference for female protagonists, likely due to overcorrection in alignment strategies, while GPT-4o achieved a more balanced representation, even incorporating non-binary protagonists. A three-way ANOVA further showed a significant interaction effect between narrative frameworks and protagonist gender on GenBit Scores, with the Heroine's Journey consistently yielding higher scores. These findings suggest that LLMs assign protagonist gender based on observed associations with specific narrative archetypes, reflecting biases embedded in training data that reflect traditional structures.

In Phase 2, the counterfactual estimation analysis exposed pronounced misclassifications when protagonist genders were swapped in narratives containing culturally gendered plot elements. For example, male protagonists in introspective or emotion-driven narratives were frequently misclassified as following the Heroine's Journey, while female protagonists in logic-driven conflicts were misclassified as Hero's Journeys. Non-binary protagonists also faced inconsistent treatment, with some counterfactual classification explanations introducing binary gender cues despite the original narrative pairs lacking explicit markers. These patterns underscore the persistence of stereotypical causal associations within LLMs, possibly originating from imbalances in training data, architectural constraints, or flawed alignment procedures.

Despite the relevant findings, some limitations must be acknowledged. The current counterfactual approach is limited to binary gender representations, making it less applicable to non-binary identities and cases where gender swapping is not appropriate, such as a pregnant protagonist. It is also important to consider that the GenBit Score is influenced by narrative characteristics, including the number of characters, which may lead to inflated scores in stories that feature fewer characters. Furthermore, the relatively small size of our factual dataset restricts the validity of the findings, while variations in narrative length and structure may introduce additional confounding factors. Addressing these challenges will require larger and more diverse LLM-generated narrative datasets, improved bias metrics, and greater transparency in the design and evaluation of generative systems.

While our findings point to distinct patterns of gender bias across different LLMs, precisely disentangling the influence of pretraining data from alignment interventions remains a methodological challenge. Given the limited public documentation of data curation and alignment objectives for most LLMs, we deliberately structured our experiments

to foreground observable effects in model outputs rather than speculate on internal mechanisms. By combining controlled narrative generation, counterfactual analysis, and statistical evaluation, our approach provides a practical framework for examining representational biases in LLM storytelling systems, even in the absence of full model transparency. Future research should build on this framework to refine evaluation methods and develop indirect strategies for auditing fairness in LLMs.

Another important direction for future research is the implementation and evaluation of gender bias mitigation strategies in LLM-based storytelling systems, such as PatternTeller [11]. While existing approaches to bias mitigation in LLMs focus on data curation and model fine-tuning [47], their effectiveness in narrative contexts remains underexplored. Given the complexity of large-scale LLMs, modifying their internal mechanisms to address bias presents significant challenges [48], including the risk of unintended alterations to language fluency, creativity, and coherence. Instead of intervening at the model level, we plan to investigate system-level strategies that introduce an intermediary layer between the storytelling system and the model. This layer could dynamically analyze and adjust prompts to encourage more balanced gender representation, ensuring diversity while maintaining narrative coherence.

Beyond its technical implications, our study highlighted how LLM-generated narratives can reinforce traditional gender norms, shaping the way characters are represented within structured frameworks. Storytelling is a powerful mechanism for defining societal perceptions, and the replication of restrictive gender norms by LLMs risks reinforcing historical imbalances rather than fostering more inclusive representations. By demonstrating how narrative structures influence gender representation in AI-driven storytelling, our findings emphasize the need for storytelling systems that are not merely reflections of existing biases but active participants in redefining narrative conventions. Addressing these challenges is not only a matter of technical refinement but a fundamental step toward ensuring that AI-generated narratives transcend traditional gender roles and representations.

To support transparency and reproducibility, the complete dataset of factual and counterfactual narratives used in this study – along with the source code for narrative generation, classification, and bias analysis – is publicly available via our Dataset Explorer webpage: <https://narrativelab.org/hero-bias-dataset/>.

## CRedit authorship contribution statement

**Irene C.E. van Blerck:** Writing – review & editing, Writing – original draft, Visualization, Validation, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Edirlei Soares de Lima:** Writing – review & editing, Writing – original draft, Validation, Supervision, Software, Project administration, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Margot M.E. Neggers:** Writing – review & editing, Validation, Methodology, Formal analysis, Conceptualization. **Toon Calders:** Writing – review & editing, Supervision.

## Declaration of competing interest

The authors declare that there are no known conflicts of interest associated with this publication and there has been no significant financial support for this work that could have influenced its outcome.

Corresponding author confirm that the manuscript has been read and approved by all named authors and that there are no other persons who satisfied the criteria for authorship but are not listed. Corresponding author further confirm that the order of authors listed in the manuscript has been approved by all of us.

Corresponding author confirm that we have given due consideration to the protection of intellectual property associated with this work and that there are no impediments to publication, including the timing of publication, with respect to intellectual property. In so doing I confirm that we have followed the regulations of our institutions concerning intellectual property.



## Appendix A. Neutral hero's journey

1. **The Ordinary World.** The protagonist is introduced in their ordinary, everyday life. They might have some desires, fears, or limitations. However, they are unaware of the adventure that awaits them.
  2. **The Call to Adventure.** The protagonist receives a call to leave their ordinary world and embark on a quest or adventure. This call can come in various forms, such as a message, a revelation, a dream, or an event that disrupts their normal life. Initially, the protagonist might refuse the call due to fear or responsibilities.
  3. **Refusal of the Call.** The protagonist hesitates and resists answering the call to adventure. They may be afraid of the unknown or reluctant to leave behind their familiar life. Often, the protagonist needs some motivation or encouragement to move forward.
  4. **Meeting the Mentor.** The protagonist encounters a mentor, a wise figure, or a guide who provides advice, training, or magical tools that will help them on their journey. The mentor gives the protagonist the confidence and knowledge they need to face the challenges ahead.
  5. **Crossing the Threshold.** After overcoming their reluctance, the protagonist finally steps out of their ordinary world and enters the special world of the adventure. This threshold signifies the protagonist's commitment to the journey and marks the beginning of significant changes in their life.
  6. **Tests, Allies, and Enemies.** In the special world, the protagonist faces a series of challenges, tests, and obstacles. They meet both allies who support their quest and enemies who oppose them. Each challenge serves as a learning experience, helping the protagonist grow and gain valuable skills.
  7. **Approach to the Inmost Cave.** The protagonist approaches the most dangerous and critical part of their journey — the inmost cave. It could be a physical location, a psychological challenge, or an inner conflict that they must confront. This stage represents a moment of great tension and self-doubt.
  8. **Ordeal.** The ordeal is the central crisis of the protagonist's journey. The protagonist faces their greatest fear, undergoes a transformation, or faces the most powerful enemy. This is a life-or-death moment that tests the protagonist's courage, determination, and growth.
  9. **Reward.** Having survived the ordeal, the protagonist emerges stronger and wiser. They receive a reward or gain new knowledge and insights. This could be a physical object, a realization, or a special power that will aid them in the final battle.
  10. **The Road Back.** After the ordeal, the protagonist begins their journey back to the ordinary world. However, the return journey is not without challenges. The protagonist may face further obstacles or temptations to stray from their path.
  11. **Resurrection.** In this climactic stage, the protagonist faces one last and most dangerous confrontation, often with the story's primary antagonist or their own inner demons. They are transformed through this experience and undergo a metaphorical death and rebirth.
  12. **Return with the Elixir.** Having overcome the final challenge, the protagonist returns to their ordinary world, bringing with them the elixir, a boon, or newfound wisdom. This may be a tangible item or something intangible that will benefit the protagonist's community or the world at large.
1. **Separation from the Inner Self.** The protagonist begins by rejecting aspects of their inner self, such as intuition, emotions, or creativity, often due to societal or cultural pressures. They prioritize external success or rational thinking, distancing themselves from their more nurturing, introspective side.
  2. **Identification with External Values and Gathering of Allies.** The protagonist adopts external values such as achievement, independence, and logic, believing that these traits will lead to success. They seek out allies who support this path, aligning with those who value external accomplishments and rationality.
  3. **Road of Trials and Gathering of Allies.** The protagonist faces challenges and obstacles on their journey toward success. As they confront these trials, they gather more allies who help them navigate the challenges, but these relationships are often based on shared values of external validation and accomplishment.
  4. **Finding the Illusory Boon of Success.** The protagonist achieves what they believed would bring fulfillment — recognition, status, or success — but soon realizes that this reward is superficial. It does not provide the deep emotional or spiritual satisfaction they were seeking, leaving them feeling unfulfilled.
  5. **Awakening to Feelings of Spiritual Death.** After realizing the emptiness of their achievements, the protagonist experiences a crisis. They feel disconnected from their true self, encountering a sense of spiritual or emotional death. This stage marks the beginning of a deeper internal journey.
  6. **Initiation and Descent into the Inner World.** The protagonist embarks on a deep inward journey, confronting aspects of themselves that they had neglected or suppressed. This could involve facing their own emotions, subconscious fears, or hidden strengths, leading to a profound understanding of their inner world.
  7. **Urgent Yearning to Reconnect with the Inner Self.** After this descent, the protagonist feels an urgent need to reconnect with their emotional and intuitive side. They seek balance between their external goals and their internal needs, striving to integrate all parts of their identity.
  8. **Healing the Split Between Past and Present.** The protagonist works to heal the rift between their past and present selves, often symbolized by reconciling with key relationships or beliefs that shaped their identity. They begin to accept and integrate the parts of themselves they once rejected or neglected.
  9. **Healing the Wounded External Self.** The protagonist addresses the wounded aspects of their external identity — traits like aggression, dominance, or control — that have become damaging. They seek to balance these traits with more nurturing, compassionate, or intuitive qualities.
  10. **Integration of Internal and External Aspects.** The protagonist achieves a balance between their external and internal aspects, no longer seeing them in opposition, and reaches a state of inner wholeness. This integration allows them to live authentically, fully embracing both their rational and intuitive sides.

## Appendix C. Example of factual narrative: Shadows of solitude

### • Chapter 1:

- Sarah, the sole survivor, wanders the deserted streets of New York City, her backpack slung over her shoulder, scavenging for food and supplies. The haunting silence of the abandoned skyscrapers looms above her, a constant reminder of the bustling city that once was.
- Each morning, Sarah wakes in the crumbling remains of what used to be her apartment, enveloped by the comforting familiarity of old photographs taped to the cracked walls, the only vestiges of a past life.

## Appendix B. Neutral heroine's journey

- Afternoons find Sarah in the New York Public Library, poring over medical texts and research notes, driven by her relentless desire to find a cure for the plague that ravaged humanity.
  - She often ventures across Central Park, careful to avoid the shadows where the transformed lurk, feeling the weight of fear tug at her as her mind battles between caution and the need for normalcy.
  - As dusk falls, Sarah lights candles in her makeshift laboratory, a space cluttered with vials and old equipment. Here, she spends her evenings experimenting, fueled by hope and longing for human connection.
- Chapter 2:
    - One fateful evening, as Sarah delves into her research, she discovers a dusty notebook hidden within the pages of an old medical journal. It contains cryptic notes and sketches detailing a potential cure for the plague.
    - Intrigued and skeptical, Sarah spends days deciphering the notebook, piecing together its fragmented information. She feels a spark of hope, tantalizingly close yet frustratingly elusive.
    - As she explores further into the journal, Sarah finds a map suggesting the existence of a hidden laboratory in the city, where the original research might have taken place. This revelation pulls at her, beckoning her to venture beyond her familiar safety.
    - That night, Sarah dreams vividly of the city alive again, free of monsters, its streets filled with laughter and light. Upon waking, she feels the compelling call to seek out the laboratory and uncover its secrets.
    - The call to adventure resonates within her, urging her to leave behind the relative safety of her routine and embrace the uncertain journey for the cure she longs to find.
  - Chapter 3:
    - Despite the allure of the notebook, Sarah hesitates, torn between the promise of discovery and the fear of venturing into the unknown. The thought of leaving her safe haven fills her with dread.
    - Memories of past encounters with the transformed plague victims flash through her mind, their monstrous faces and haunting cries serving as a stark reminder of the dangers that lurk outside.
    - She convinces herself that her research in the safety of her apartment is enough, trying to quell the feeling of inadequacy and her growing doubts about the path that lies ahead.
    - As days pass, Sarah confines herself to her routine, grappling with her reluctance to leave, yet unable to shake the nagging feeling that she is on the brink of something momentous.
    - One evening, the sight of a distant, flickering light across the city skyline reignites a smoldering curiosity in Sarah, hinting at the possibility that she is not as alone as she believes.
  - Chapter 4:
    - One stormy night, Sarah is jolted awake by the sound of rustling papers. She finds an old radio crackling to life, broadcasting a garbled voice offering guidance and hope to any survivors seeking a cure.
    - Intrigued by the mysterious voice, Sarah listens intently, the voice speaking of scientific knowledge and strategies for evading the transformed, giving her newfound courage and inspiration.
  - Chapter 5:
    - Through the radio, the mentor reveals stories of resilience and survival, sharing tales of those who braved dangers and made breakthroughs, strengthening Sarah's resolve to pursue the quest.
    - The voice provides cryptic hints about the location of the hidden laboratory, weaving a tapestry of clues that Sarah must unravel, guiding her through the process of mapping her path.
    - Empowered by the mentor's guidance, Sarah gathers her courage, feeling a sense of connection with the unseen voice, as if it has become an invisible ally on her journey.
  - Chapter 6:
    - With the first light of dawn breaking, Sarah takes her first steps beyond the familiar confines of her apartment, leaving behind the safety she has clung to for so long and embracing the journey ahead.
    - She navigates through the desolate streets with heightened awareness, each step crossing her further from her ordinary world and deeper into the adventure that awaits.
    - Along the way, she encounters remnants of the past – abandoned cars, forgotten belongings – which serve as poignant reminders of the world that once thrived.
    - As she approaches the threshold of the unknown, Sarah pauses at a barricade erected during the early days of the plague, a symbolic boundary she must cross to truly enter the adventure.
    - Summoning her resolve and recalling the voice of the mentor, she climbs over the barricade, leaving behind the life she knew and stepping fully into the uncharted territory of her quest.
  - Chapter 7:
    - As Sarah navigates the labyrinthine streets, her first challenge arises in the form of a group of transformed creatures blocking her path, forcing her to employ stealth and quick thinking to evade them.
    - Venturing deeper into the city, she meets Marcus, another survivor who has made a life scavenging among the ruins. Initially wary, they slowly build trust, each finding solace in the other's presence.
    - Together, they face hostile gangs of survivors who have succumbed to desperation and violence, barely escaping ambush thanks to Marcus's knowledge of hidden passages and safe zones.
    - In the quieter moments, Sarah and Marcus share stories of their past lives, deepening their bond and reigniting their hope, while forging plans to find and reach the fabled laboratory.
    - Along their journey, they cross paths with a mysterious old woman who offers cryptic advice and a small vial containing a substance she claims can protect them from the plague, leaving them more questions than answers.
  - Chapter 8:
    - Armed with the map and the old woman's vial, Sarah and Marcus make their way toward the rumored entrance to the hidden laboratory, their path fraught with obstacles as they delve deeper into the city's heart.
    - They find themselves at the entrance to the subway tunnels, which they must traverse to reach the laboratory, a daunting labyrinth of darkness and uncertainty stretching out before them.

- As they descend into the dimly lit tunnels, their senses sharpen, every sound amplified in the oppressive silence, forcing them to rely on each other as they navigate this treacherous path.
- They encounter unexpected allies — a small group of survivors living underground who provide crucial guidance on the safest routes and warn of lurking dangers, reinforcing Sarah and Marcus's resolve.
- As they near the inmost cave, the tension escalates with the realization that they are closer than ever to their goal, yet aware that the greatest challenges still lie ahead.

• Chapter 8:

- Deep within the labyrinthine subway tunnels, as Sarah and Marcus near the laboratory's rumored location, they are ambushed by a horde of transformed creatures, forcing them into a desperate battle for survival.
- Sarah's fear is palpable, but the old woman's vial provides unexpected assistance — a vapor that repels the creatures momentarily, giving them a crucial chance to escape the immediate threat.
- The reprieve is short-lived as the creatures regroup, and Sarah and Marcus must navigate a series of collapsing tunnels, each step bringing them closer to exhaustion and peril.
- They reach a dead-end, forced into a final confrontation with the creatures. Outnumbered and trapped, they muster every ounce of strength and resourcefulness to fight their way through.
- Just as they reach the brink of defeat, a hidden door is revealed by a fortuitous burst of light, leading them into the sanctuary of the laboratory, safe for the moment but shaken by their ordeal.

• Chapter 9:

- Inside the hidden laboratory, Sarah and Marcus find shelves lined with research materials and equipment left untouched since the city's fall, a treasure trove of information and potential breakthroughs in the fight against the plague.
- As Sarah examines the notebooks, she discovers the missing pieces to the cryptic notes she found earlier. They reveal a formula that could potentially lead to a cure, filling her with a renewed sense of purpose.
- Marcus finds a cache of preserved food and medical supplies, offering immediate relief and ensuring their survival as they work diligently to synthesize the potential cure.
- Working together, Sarah and Marcus manage to produce a small batch of the experimental cure, each vial representing hope for the future and the possibility of reclaiming their world.
- With the formula complete and the first batch of the cure in hand, Sarah and Marcus feel a profound sense of achievement and the knowledge that they hold the key to saving countless lives.

• Chapter 10:

- With the vials of the experimental cure secured, Sarah and Marcus prepare to leave the sanctuary of the laboratory, knowing they must return to the surface and share their discovery with those who remain.
- Retracing their steps through the subway tunnels, they move cautiously, aware that the transformed creatures and other dangers still lurk in the shadows, their senses heightened by the burden of their mission.

- Emerging from the labyrinth of tunnels, they find the cityscape eerily quiet, the atmosphere tense with a foreboding sense of the challenges yet to come as they seek to reach the survivor enclave.
- Along the way, they encounter familiar faces — other survivors whose stories they crossed — offering solace and encouragement, reminding Sarah and Marcus of the vital importance of their mission.
- As they near their destination, the weight of their responsibility grows heavier, each step on the road back filled with both the anticipation of what their discovery could mean and the fear of what might hinder their path.

• Chapter 11:

- As Sarah and Marcus near the survivor enclave, they are ambushed by a rogue faction that seeks to take the cure for themselves, forcing them into a tense standoff that tests their resolve and ingenuity.
- In the midst of the chaos, Sarah is separated from Marcus, faced with her deepest fears as she struggles to protect the vials and herself, realizing this is her moment of rebirth and true courage.
- Drawing on the skills and strength she has acquired throughout her journey, Sarah outsmarts their attackers with a carefully laid trap, managing to secure the safety of the cure and reunite with Marcus.
- This trial by fire fortifies their bond, each realizing the depths of their resilience and the power of their partnership, emerging from the confrontation stronger and more united than before.
- With the rogue threat neutralized, Sarah and Marcus press on to the enclave, carrying not just the cure, but a profound transformation within themselves, ready to usher in a new era of hope.

• Chapter 12:

- Sarah and Marcus arrive at the survivor enclave, greeted by a hopeful gathering of weary faces. They present the vials of the experimental cure, sharing the promise of a new beginning with those who have endured so much.
- The enclave's healers begin administering the cure to ailing residents, and slowly, signs of recovery emerge, confirming the efficacy of the formula and igniting a wave of joyous disbelief.
- As the days pass, the enclave transforms; laughter and renewed energy permeate the once somber atmosphere, each recovered person a testament to Sarah and Marcus's arduous journey.
- With their mission accomplished, Sarah and Marcus find themselves embraced by the community, hailed as heroes, yet humbled by the resilience and courage they have witnessed in others.
- As they stand on the verge of a new chapter, Sarah reflects on her transformation and the power of hope, knowing that she and Marcus have helped change the course of life for many, and that their story has only just begun.

#### Appendix D. Example of counterfactual narrative: Shadows of solitude

• Chapter 1:

- Sam, the sole survivor, wanders the deserted streets of New York City, his backpack slung over his shoulder, scavenging for food and supplies. The haunting silence of the abandoned skyscrapers looms above him, a constant reminder of the bustling city that once was.
  - Each morning, Sam wakes in the crumbling remains of what used to be his apartment, enveloped by the comforting familiarity of old photographs taped to the cracked walls, the only vestiges of a past life.
  - Afternoons find Sam in the New York Public Library, poring over medical texts and research notes, driven by his relentless desire to find a cure for the plague that ravaged humanity.
  - He often ventures across Central Park, careful to avoid the shadows where the transformed lurk, feeling the weight of fear tug at him as his mind battles between caution and the need for normalcy.
  - As dusk falls, Sam lights candles in his makeshift laboratory, a space cluttered with vials and old equipment. Here, he spends his evenings experimenting, fueled by hope and longing for human connection.
- Chapter 2:
    - One fateful evening, as Sam delves into his research, he discovers a dusty notebook hidden within the pages of an old medical journal. It contains cryptic notes and sketches detailing a potential cure for the plague.
    - Intrigued and skeptical, Sam spends days deciphering the notebook, piecing together its fragmented information. He feels a spark of hope, tantalizingly close yet frustratingly elusive.
    - As he explores further into the journal, Sam finds a map suggesting the existence of a hidden laboratory in the city, where the original research might have taken place. This revelation pulls at him, beckoning him to venture beyond his familiar safety.
    - That night, Sam dreams vividly of the city alive again, free of monsters, its streets filled with laughter and light. Upon waking, he feels the compelling call to seek out the laboratory and uncover its secrets.
    - The call to adventure resonates within him, urging him to leave behind the relative safety of his routine and embrace the uncertain journey for the cure he longs to find.
  - Chapter 3:
    - Despite the allure of the notebook, Sam hesitates, torn between the promise of discovery and the fear of venturing into the unknown. The thought of leaving his safe haven fills him with dread.
    - Memories of past encounters with the transformed plague victims flash through his mind, their monstrous faces and haunting cries serving as a stark reminder of the dangers that lurk outside.
    - He convinces himself that his research in the safety of his apartment is enough, trying to quell the feeling of inadequacy and his growing doubts about the path that lies ahead.
    - As days pass, Sam confines himself to his routine, grappling with his reluctance to leave, yet unable to shake the nagging feeling that he is on the brink of something momentous.
    - One evening, the sight of a distant, flickering light across the city skyline reignites a smoldering curiosity in Sam, hinting at the possibility that he is not as alone as he believes.
  - Chapter 4:
    - One stormy night, Sam is jolted awake by the sound of rustling papers. He finds an old radio crackling to life, broadcasting a garbled voice offering guidance and hope to any survivors seeking a cure.
    - Intrigued by the mysterious voice, Sam listens intently, the voice speaking of scientific knowledge and strategies for evading the transformed, giving him newfound courage and inspiration.
    - Through the radio, the mentor reveals stories of resilience and survival, sharing tales of those who braved dangers and made breakthroughs, strengthening Sam's resolve to pursue the quest.
    - The voice provides cryptic hints about the location of the hidden laboratory, weaving a tapestry of clues that Sam must unravel, guiding him through the process of mapping his path.
    - Empowered by the mentor's guidance, Sam gathers his courage, feeling a sense of connection with the unseen voice, as if it has become an invisible ally on his journey.
  - Chapter 5:
    - With the first light of dawn breaking, Sam takes his first steps beyond the familiar confines of his apartment, leaving behind the safety he has clung to for so long and embracing the journey ahead.
    - He navigates through the desolate streets with heightened awareness, each step crossing him further from his ordinary world and deeper into the adventure that awaits.
    - Along the way, he encounters remnants of the past – abandoned cars, forgotten belongings – which serve as poignant reminders of the world that once thrived.
    - As he approaches the threshold of the unknown, Sam pauses at a barricade erected during the early days of the plague, a symbolic boundary he must cross to truly enter the adventure.
    - Summoning his resolve and recalling the voice of the mentor, he climbs over the barricade, leaving behind the life he knew and stepping fully into the uncharted territory of his quest.
  - Chapter 6:
    - As Sam navigates the labyrinthine streets, his first challenge arises in the form of a group of transformed creatures blocking his path, forcing him to employ stealth and quick thinking to evade them.
    - Venturing deeper into the city, he meets Marcia, another survivor who has made a life scavenging among the ruins. Initially wary, they slowly build trust, each finding solace in the other's presence.
    - Together, they face hostile gangs of survivors who have succumbed to desperation and violence, barely escaping ambush thanks to Marcia's knowledge of hidden passages and safe zones.
    - In the quieter moments, Sam and Marcia share stories of their past lives, deepening their bond and reigniting their hope, while forging plans to find and reach the fabled laboratory.
    - Along their journey, they cross paths with a mysterious old



man who offers cryptic advice and a small vial containing a substance he claims can protect them from the plague, leaving them more questions than answers.

• Chapter 7:

- Armed with the map and the old man's vial, Sam and Marcia make their way toward the rumored entrance to the hidden laboratory, their path fraught with obstacles as they delve deeper into the city's heart.
- They find themselves at the entrance to the subway tunnels, which they must traverse to reach the laboratory, a daunting labyrinth of darkness and uncertainty stretching out before them.
- As they descend into the dimly lit tunnels, their senses sharpen, every sound amplified in the oppressive silence, forcing them to rely on each other as they navigate this treacherous path.
- They encounter unexpected allies — a small group of survivors living underground who provide crucial guidance on the safest routes and warn of lurking dangers, reinforcing Sam and Marcia's resolve.
- As they near the inmost cave, the tension escalates with the realization that they are closer than ever to their goal, yet aware that the greatest challenges still lie ahead.

• Chapter 8:

- Deep within the labyrinthine subway tunnels, as Sam and Marcia near the laboratory's rumored location, they are ambushed by a horde of transformed creatures, forcing them into a desperate battle for survival.
- Sam's fear is palpable, but the old man's vial provides unexpected assistance — a vapor that repels the creatures momentarily, giving them a crucial chance to escape the immediate threat.
- The reprieve is short-lived as the creatures regroup, and Sam and Marcia must navigate a series of collapsing tunnels, each step bringing them closer to exhaustion and peril.
- They reach a dead-end, forced into a final confrontation with the creatures. Outnumbered and trapped, they muster every ounce of strength and resourcefulness to fight their way through.
- Just as they reach the brink of defeat, a hidden door is revealed by a fortuitous burst of light, leading them into the sanctuary of the laboratory, safe for the moment but shaken by their ordeal.

• Chapter 9:

- Inside the hidden laboratory, Sam and Marcia find shelves lined with research materials and equipment left untouched since the city's fall, a treasure trove of information and potential breakthroughs in the fight against the plague.
- As Sam examines the notebooks, he discovers the missing pieces to the cryptic notes he found earlier. They reveal a formula that could potentially lead to a cure, filling him with a renewed sense of purpose.
- Marcia finds a cache of preserved food and medical supplies, offering immediate relief and ensuring their survival as they work diligently to synthesize the potential cure.
- Working together, Sam and Marcia manage to produce a small batch of the experimental cure, each vial representing hope for the future and the possibility of reclaiming their world.
- With the formula complete and the first batch of the cure in hand, Sam and Marcia feel a profound sense of achievement and the knowledge that they hold the key to saving countless lives.

• Chapter 10:

- With the vials of the experimental cure secured, Sam and Marcia prepare to leave the sanctuary of the laboratory, knowing they must return to the surface and share their discovery with those who remain.
- Retracing their steps through the subway tunnels, they move cautiously, aware that the transformed creatures and other dangers still lurk in the shadows, their senses heightened by the burden of their mission.
- Emerging from the labyrinth of tunnels, they find the cityscape eerily quiet, the atmosphere tense with a foreboding sense of the challenges yet to come as they seek to reach the survivor enclave.
- Along the way, they encounter familiar faces — other survivors whose stories they crossed — offering solace and encouragement, reminding Sam and Marcia of the vital importance of their mission.
- As they near their destination, the weight of their responsibility grows heavier, each step on the road back filled with both the anticipation of what their discovery could mean and the fear of what might hinder their path.

• Chapter 11:

- As Sam and Marcia near the survivor enclave, they are ambushed by a rogue faction that seeks to take the cure for themselves, forcing them into a tense standoff that tests their resolve and ingenuity.
- In the midst of the chaos, Sam is separated from Marcia, faced with his deepest fears as he struggles to protect the vials and himself, realizing this is his moment of rebirth and true courage.
- Drawing on the skills and strength he has acquired throughout his journey, Sam outsmarts their attackers with a carefully laid trap, managing to secure the safety of the cure and reunite with Marcia.
- This trial by fire fortifies their bond, each realizing the depths of their resilience and the power of their partnership, emerging from the confrontation stronger and more united than before.
- With the rogue threat neutralized, Sam and Marcia press on to the enclave, carrying not just the cure, but a profound transformation within themselves, ready to usher in a new era of hope.

• Chapter 12:

- Sam and Marcia arrive at the survivor enclave, greeted by a hopeful gathering of weary faces. They present the vials of the experimental cure, sharing the promise of a new beginning with those who have endured so much.
- The enclave's healers begin administering the cure to ailing residents, and slowly, signs of recovery emerge, confirming the efficacy of the formula and igniting a wave of joyous disbelief.
- As the days pass, the enclave transforms; laughter and renewed energy permeate the once somber atmosphere, each recovered person a testament to Sam and Marcia's arduous journey.
- With their mission accomplished, Sam and Marcia find themselves embraced by the community, hailed as heroes, yet humbled by the resilience and courage they have witnessed in others.
- As they stand on the verge of a new chapter, Sam reflects on his transformation and the power of hope, knowing that he and Marcia have helped change the course of life for many, and that their story has only just begun.

## Data availability

Data will be made available on request.

## References

- [1] M. Imran, N. Almusharraf, Analyzing the role of ChatGPT as a writing assistant at higher education level: A systematic review of the literature, *Contemp. Educ. Technol.* 15 (4) (2023) ep464, <http://dx.doi.org/10.30935/cedtech/13605>.
- [2] J. Regalia, From briefs to bytes: How generative AI is transforming legal writing and practice, *Tulsa Law Rev.* 59 (2024) URL <https://digitalcommons.law.utulsa.edu/tlr/vol59/iss2/5>.
- [3] A. Calderwood, V. Qiu, K. Gero, L.B. Chilton, How novelists use generative language models: An exploratory user study, in: *Joint Proceedings of the Workshops on Human-AI Co-Creation with Generative Models and User-Aware Conversational Agents*, 25th International Conference on Intelligent User Interfaces, IUI 2020, 2020, pp. 1–5, URL <https://ceur-ws.org/Vol-2848/HAI-GEN-Paper-3.pdf>.
- [4] D. Ippolito, A. Yuan, A. Coenen, S. Burnam, Creative writing with an AI-powered writing assistant: Perspectives from professional writers, in: *ArXiv E-Prints*, ArXiv: 2211.05030 [Cs.CL], 2022, pp. 1–18, <http://dx.doi.org/10.48550/arXiv.2211.05030>, arXiv:2211.05030.
- [5] H. Al Naqbi, Z. Bahroun, V. Ahmed, Enhancing work productivity through generative artificial intelligence: A comprehensive literature review, *Sustainability* 16 (3) (2024) <http://dx.doi.org/10.3390/su16031166>.
- [6] G. Franceschelli, M. Musolesi, On the creativity of large language models, *AI & SOCIETY* (2024) <http://dx.doi.org/10.1007/s00146-024-02127-3>.
- [7] L. Castricato, S. Frazier, J. Balloch, M. Riedl, Tell me a story like i'm five: Story generation via question answering, in: *Proceedings of the 3rd Workshop on Narrative Understanding*, 2021, pp. 1–8, URL <https://par.nsf.gov/biblio/10249509>.
- [8] E.S. de Lima, M.A. Casanova, A.L. Furtado, Imagining from images with an AI storytelling tool, in: *ArXiv E-Prints*, ArXiv: 2408.11517 [Cs.CL], 2024, pp. 1–16, <http://dx.doi.org/10.48550/arXiv.2408.11517>, arXiv:2408.11517.
- [9] E.S. de Lima, B. Feijó, M.A. Casanova, A.L. Furtado, ChatGeppetto - an AI-powered storyteller, in: *Proceedings of the 22nd Brazilian Symposium on Games and Digital Entertainment*, ACM, 2024, pp. 28–37, <http://dx.doi.org/10.1145/3631085.3631302>.
- [10] E.S. de Lima, M.M. Neggers, B. Feijó, M.A. Casanova, A.L. Furtado, An AI-powered approach to the semiotic reconstruction of narratives, *Entertain. Comput.* 52 (2025) 100810, <http://dx.doi.org/10.1016/j.entcom.2024.100810>.
- [11] E.S. de Lima, M.M.E. Neggers, M.A. Casanova, B. Feijó, A.L. Furtado, A pattern-oriented AI-powered approach to story composition, in: P. Figueroa, et al. (Eds.), *Entertainment Computing – ICEC 2024*, Springer Cham, 2024, pp. 1–16, [http://dx.doi.org/10.1007/978-3-031-74353-5\\_11](http://dx.doi.org/10.1007/978-3-031-74353-5_11).
- [12] S. Värtinen, P. Hämäläinen, C. Guckelsberger, Generating role-playing game quests with GPT language models, *IEEE Trans. Games* 16 (1) (2024) 127–139, <http://dx.doi.org/10.1109/TG.2022.3228480>.
- [13] W. Babonnaud, E. Delouche, M. Lahlouh, The bias that lies beneath: Qualitative uncovering of stereotypes in large language models, in: *14th Scandinavian Conference on Artificial Intelligence SCAI 2024*, 2024, pp. 195–203, <http://dx.doi.org/10.3384/ecp208022>.
- [14] R. Navigli, S. Conia, B. Ross, Biases in large language models: Origins, inventory, and discussion, *J. Data Inf. Qual.* 15 (2) (2023) 1–21, <http://dx.doi.org/10.1145/3597307>.
- [15] E. Ntoutsis, P. Fafalios, U. Gadiraju, V. Iosifidis, W. Nejdl, M.-E. Vidal, S. Ruggieri, F. Turini, S. Papadopoulos, E. Krasanakis, et al., Bias in data-driven artificial intelligence systems—An introductory survey, *Wiley Interdiscip. Rev.: Data Min. Knowl. Discov.* 10 (3) (2020) e1356, <http://dx.doi.org/10.1002/widm.1356>.
- [16] E.S. de Lima, B. Feijó, A.L. Furtado, Managing the plot structure of character-based interactive narratives in games, *Entertain. Comput.* 47 (2023) 100590, <http://dx.doi.org/10.1016/j.entcom.2023.100590>.
- [17] J. Campbell, *The Hero with a Thousand Faces*, New World Library, 2008.
- [18] M. Murdock, *The Heroine's Journey: Woman's Quest for Wholeness*, Shambhala, 1990.
- [19] B. Ip, Narrative structures in computer and video games: Part 1: Context, definitions, and initial findings, *Games Cult.* 6 (2) (2011) 103–134, <http://dx.doi.org/10.1177/1555412010364982>.
- [20] G. Carriger, *The Heroine's Journey: for Writers, Readers, and Fans of Pop Culture*, Gail Carriger, 2020, eBook.
- [21] L. Lucy, D. Bamman, Gender and representation bias in GPT-3 generated stories, in: *Proceedings of the Third Workshop on Narrative Understanding*, Association for Computational Linguistics, 2021, pp. 48–55, <http://dx.doi.org/10.18653/v1/2021.nuse-1.5>.
- [22] D. Jackson, M. Courneya, Unreliable narrator: Reparative approaches to harmful biases in AI storytelling for the HE classroom and future creative industries, *Braz. Creative Ind. J.* 3 (2) (2023) 59–75, <http://dx.doi.org/10.25112/bcij.v3i2.3540>.
- [23] N. Beguš, Experimental narratives: A comparison of human crowdsourced storytelling and AI storytelling, *Humanit. Soc. Sci. Commun.* 11 (1) (2024) 1392, <http://dx.doi.org/10.1057/s41599-024-03868-8>.
- [24] J.H. Miller, *Versions of Pygmalion*, Harvard University Press, 1990.
- [25] A. Marin, M. Eger, Towards evaluating profession-based gender bias in ChatGPT and its impact on narrative generation, in: M. Farrokhmaleki, P. Rahmati, K. Saadat, R. Zhao (Eds.), *Proceedings of the AIIDE Workshop on Intelligent Narrative Technologies Co-located with the 20th AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, AIIDE 2024, Lexington, Kentucky, USA, 2024, pp. 1–8.
- [26] P. Taveekitworachai, F. Abdullah, M.C. Gursesli, M.F. Dewantoro, S. Chen, A. Lanata, A. Guazzini, R. Thawonmas, What is waiting for us at the end? Inherent biases of game story endings in large language models, in: L. Holloway-Attaway, J.T. Murray (Eds.), *Interactive Storytelling*, Springer Nature Switzerland, Cham, 2023, pp. 274–284, [http://dx.doi.org/10.1007/978-3-031-47658-7\\_26](http://dx.doi.org/10.1007/978-3-031-47658-7_26).
- [27] P. Taveekitworachai, K. Plupattanakit, R. Thawonmas, Assessing inherent biases following prompt compression of large language models for game story generation, in: *2024 IEEE Conference on Games, CoG*, 2024, pp. 1–4, <http://dx.doi.org/10.1109/CoG60054.2024.10645609>.
- [28] I. Weissburg, S. Anand, S. Levy, H. Jeong, LLMs are biased teachers: Evaluating LLM bias in personalized education, in: *ArXiv E-Prints*, ArXiv: 2410.14012 [Cs.CL], 2024, pp. 1–49, <http://dx.doi.org/10.48550/arXiv.2410.14012>, arXiv:2410.14012.
- [29] A. Kwako, C. Ormerod, Can language models guess your identity? Analyzing demographic biases in AI essay scoring, in: E. Kochmar, M. Bexte, J. Burstein, A. Horbach, R. Laarmann-Quante, A. Tack, V. Yaneva, Z. Yuan (Eds.), *Proceedings of the 19th Workshop on Innovative Use of NLP for Building Educational Applications*, BEA 2024, Association for Computational Linguistics, Mexico City, Mexico, 2024, pp. 78–86, URL <https://aclanthology.org/2024.bea-1.7/>.
- [30] Y. Huang, Unveiling gender bias in large language models: Using teacher's evaluation in higher education as an example, in: *ArXiv E-Prints*, ArXiv: 2409.09652 [Cs.CL], 2024, pp. 1–30, <http://dx.doi.org/10.48550/arXiv.2409.09652>, arXiv:2409.09652.
- [31] L. Lippens, Computer says 'no': Exploring systemic bias in ChatGPT using an audit approach, *Comput. Hum. Behav.: Artif. Hum.* 2 (1) (2024) 100054, <http://dx.doi.org/10.1016/j.chbah.2024.100054>.
- [32] L. Armstrong, A. Liu, S. MacNeil, D. Metaxa, The silicon ceiling: Auditing GPT's race and gender biases in hiring, in: *Proceedings of the 4th ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO '24, Association for Computing Machinery, New York, NY, USA, 2024, pp. 1–18, <http://dx.doi.org/10.1145/3689904.3694699>.
- [33] T. Zack, E. Lehman, M. Suzgun, J.A. Rodriguez, L.A. Celi, J. Gichoya, D. Jurafsky, P. Szolovits, D.W. Bates, R.-E.E. Abdunour, A.J. Butte, E. Alsentzer, Assessing the potential of GPT-4 to perpetuate racial and gender biases in health care: a model evaluation study, *Lancet Digit. Heal.* 6 (1) (2024) e12–e22, [http://dx.doi.org/10.1016/S2589-7500\(23\)00225-X](http://dx.doi.org/10.1016/S2589-7500(23)00225-X).
- [34] S. Rickman, Evaluating gender bias in large language models in long-term care, *Prepr. Available At Res. Sq.* (2024) <http://dx.doi.org/10.21203/rs.3.rs-5166499/v2>.
- [35] B.D. Menz, N.M. Kuderer, B. Chin-Yee, J.M. Logan, A. Rowland, M.J. Sorich, A.M. Hopkins, Gender representation of health care professionals in large language model-generated stories, *JAMA Netw. Open* 7 (9) (2024) e2434997, <http://dx.doi.org/10.1001/jamanetworkopen.2024.34997>.
- [36] A. Agrawal, Fairness in AI-driven oncology: Investigating racial and gender biases in large language models, *Cureus* 16 (9) (2024) e69541, <http://dx.doi.org/10.7759/cureus.69541>.
- [37] E.S. de Lima, M.M.E. Neggers, A.L. Furtado, Multigenre AI-powered story composition, in: *ArXiv e-prints*, ArXiv: 2405.06685 [Cs.CL], 2024, pp. 1–17, <http://dx.doi.org/10.48550/arXiv.2405.06685>, arXiv:2405.06685.
- [38] S. Wang, G. Durrett, K. Erk, Narrative interpolation for generating and understanding stories, in: *ArXiv E-Prints*, ArXiv: 2008.07466 [Cs.CL], 2020, pp. 1–5, <http://dx.doi.org/10.48550/arXiv.2008.07466>, arXiv:2008.07466.
- [39] T. Frontgia, Archetypes, stereotypes, and the female hero: Transformations in contemporary perspectives, *Mythlore* 18 (1991) 15–18, URL <http://www.jstor.org/stable/26812483>.
- [40] M. Murdock, *The Heroine's Journey Workbook: A Map for Every Woman's Quest*, Shambhala, 2020.
- [41] K. Sengupta, R. Maher, D. Groves, C. Olieman, GenBiT: measure and mitigate gender bias in language datasets, *Microsoft J. Appl. Res.* 16 (2021) 63–71.
- [42] D.M. Ziegler, N. Stiennon, J. Wu, T.B. Brown, A. Radford, D. Amodei, P. Christiano, G. Irving, Fine-tuning language models from human preferences, in: *ArXiv E-Prints*, ArXiv: 1909.08593 [Cs.CL], 2020, pp. 1–26, <http://dx.doi.org/10.48550/arXiv.1909.08593>, arXiv:1909.08593.
- [43] H. Touvron, et al., Llama 2: Open foundation and fine-tuned chat models, in: *ArXiv E-Prints*, ArXiv: 2307.09288 [Cs.CL], 2023, pp. 1–77, <http://dx.doi.org/10.48550/arXiv.2307.09288>, arXiv:2307.09288.

- [44] A. Grattafiori, et al., The llama 3 herd of models, in: ArXiv E-Prints, ArXiv: 2407.21783 [Cs.AI], 2024, pp. 1–92, <http://dx.doi.org/10.48550/arXiv.2407.21783>, arXiv:2407.21783.
- [45] K.K. Chang, M. Cramer, S. Soni, D. Bamman, Speak, memory: An archaeology of books known to ChatGPT/GPT-4, in: ArXiv E-Prints, ArXiv: 2305.00118 [Cs.CL], 2023, pp. 1–16, <http://dx.doi.org/10.48550/arXiv.2305.00118>, arXiv: 2305.00118.
- [46] M. Foucault, *The History of Sexuality, Volume 1: An Introduction*, Pantheon Books, New York, 1978, Translated by Robert Hurley.
- [47] I.O. Gallegos, R.A. Rossi, J. Barrow, M.M. Tanjim, S. Kim, F. Dernoncourt, T. Yu, R. Zhang, N.K. Ahmed, Bias and fairness in large language models: A survey, *Comput. Linguist.* 50 (3) (2024) 1097–1179, [http://dx.doi.org/10.1162/coli\\_a\\_00524](http://dx.doi.org/10.1162/coli_a_00524).
- [48] Z. Lin, S. Guan, W. Zhang, H. Zhang, Y. Li, H. Zhang, Towards trustworthy LLMs: a review on debiasing and dehallucinating in large language models, *Artif. Intell. Rev.* 57 (9) (2024) 243, <http://dx.doi.org/10.1007/s10462-024-10896-y>.