

Do Large Language Models Understand Mansplaining? *Well, actually...*

Carla Perez-Almendros and Jose Camacho-Collados

School of Computer Science & Informatics
Cardiff University, UK
{perezalmendrosc, camachocollados}@cardiff.ac.uk

Abstract

Gender bias has been widely studied by the NLP community. However, other more subtle variations of it, such as mansplaining, have yet received little attention. Mansplaining is a discriminatory behaviour that consists of a condescending treatment or discourse towards women. In this paper, we introduce and analyze *Well, actually...*, a corpus of 886 mansplaining stories experienced by women. We analyze the corpus in terms of features such as offensiveness, sentiment or misogyny, among others. We also explore to what extent Large Language Models (LLMs) can understand and identify mansplaining and other gender-related microaggressions. Specifically, we experiment with ChatGPT-3.5-Turbo and LLaMA-2 (13b and 70b), with both targeted and open questions. Our findings suggest that, although they can identify mansplaining to some extent, LLMs still struggle to point out this attitude and will even reproduce some of the social patterns behind mansplaining situations, for instance by praising men for giving unsolicited advice to women.

Keywords: mansplaining, corpus linguistics, large language models, social media, gender bias

1. Introduction

Since the introduction of transformers and pre-trained Language Models (LMs) (Vaswani et al., 2017), researchers in NLP have worked on identifying and mitigating different types of human biases (Lauscher and Glavaš, 2019; Caliskan et al., 2017; Bolukbasi et al., 2016; Manzini et al., 2019). These works have focused on the biases contained in training data, as well as on the biases reproduced or even amplified by the models themselves (Zhao et al., 2018; Hovy and Prabhumoye, 2021; Wang and Russakovsky, 2021; Hall et al., 2022; Zhou et al., 2023). With the introduction of Large Language Models (LLMs), and their growing involvement in our daily life, addressing these biases seems more important than ever. On the one hand, LLMs, as well as their predecessors, contain biases, as they are trained on vast amounts of (biased) data online (Abid et al., 2021; Schramowski et al., 2022; Gallegos et al., 2023). On the other hand, there is an evident effort by the research community in identifying these biases and in trying to mitigate them in order to obtain better and fairer models (Navigli et al., 2023; Steed et al., 2022; Ferrara, 2023). As a case in point, evaluating and mitigating gender biases has been an active area of research in many NLP tasks, such as text classification, machine translation or text generation (Sheng et al., 2019; Savoldi et al., 2021; Kotek et al., 2023). However, there may be some types of biases which are more difficult to detect and mitigate than others, as they are also for humans. For instance, NLP models might struggle to identify biases based on a more subtle language or those without a direct intention to harm others (Or-

tiz, 2022; Perez-Almendros and Schockaert, 2022; Naous et al., 2023; Kiehne et al., 2022; Ramezani and Xu, 2023).

In this paper, we focus on mansplaining, a type of (subtle) discriminatory behaviour based on a biased attitude towards women (Dular, 2021; Briggs et al., 2023). We can find a situation of mansplaining when a man treats or explains something to a woman in a manner regarded as condescending or patronizing. This treatment comes motivated by the assumption that the woman, just because of her gender, will not have the knowledge about the topic or the capacity to understand it without the help of a man, or the belief that the woman cannot have more expertise in the topic than a man. With this work, we introduce two main contributions to the community.

First, we collect, curate and release "*Well, actually...*", a corpus of 886 short stories about mansplaining situations. We perform a qualitative analysis of these stories with the objective of better understanding the language involved when describing mansplaining situations, as well as other types of biases which might support mansplaining attitudes. Specifically, we analyze the data to detect offensive language, sentiments, irony, condescension and misogyny. We also explore the most popular topics in the data to better understand where we can encounter mansplaining.

Second, we analyze to what extent LLMs understand mansplaining. Specifically, we ask Llama-2-13b, Llama-2-70b and ChatGPT-3.5-Turbo (hereafter referred to as ChatGPT) several questions to identify whether they are able to detect gender-related discriminatory attitudes in the

stories contained in the *Well, actually...* corpus. Furthermore, we manually annotate some of the models' answers for comparison and perform qualitative analysis over some examples. Overall, we find that 1) models which target other types of biased behaviour are not sufficient to model the language involved in mansplaining stories, and 2) LLMs can to some extent identify gender discrimination if explicitly asked to look for it, but they fail to detect it if asked more general questions.

2. Related Work

NLP has been targeting human biases for a long time, and tasks such as sentiment analysis (Patwa et al., 2020; Barnes et al., 2022), hate speech (Basile et al., 2019), offensive language (Zampieri et al., 2020) or gender bias detection (Sun et al., 2019; Fersini et al., 2022; Kotek et al., 2023) have become well-established challenges in the community.

However, other types of equally harmful but more subtle discriminatory language still pose a challenge for the NLP community. One of such cases is microaggressions, where the discriminatory treatment is usually subtle and indirect. They tend to reproduce biases deeply accepted by society and their identification is usually subjective (Breitfeller et al., 2019; Ògúnremí et al., 2022). In their work about "Finding Microaggressions in the Wild", Breitfeller et al. (2019) create a taxonomy of types of microaggressions, with some of them, such as *stereotypes*, *objectification*, *second-class citizen* or *abnormality*, specially affecting women. The discourse of condescension, which also relies in subtle and subjective language, often concealed of good intentions, has also been recently studied by the NLP community (Perez-Almendros et al., 2022). Works in this area have highlighted how NLP systems still struggle to identify condescension and the necessity to focus on such subtle types of discriminatory treatments (Wang and Potts, 2019; Perez-Almendros and Schockaert, 2022).

Mansplaining can be considered as a type of microaggression which uses a condescending and patronizing discourse or treatment towards women (Lutzky and Lawson, 2019; Eschmann et al., 2020; Dular, 2021). Researchers such as Lutzky and Lawson (2019) and Koc-Michalska et al. (2021) have studied this phenomenon from a data analysis and computational linguistics point of view. However, to the best of our knowledge, the NLP community has still not focused on analyzing how current NLP systems, and in particular LLMs, understand this type of discriminatory attitude. Moreover, there is no established benchmark or dataset specifically focused on mansplaining situations openly available. In this paper, we try to fill this gap by both compiling

a benchmark of mansplaining stories and analysing it, particularly from the perspective of LLMs.

3. The *Well, actually...* Corpus

The *Well, actually...* corpus contains 886 short stories, where the protagonist, usually self-identified as a female, explains their own or someone else's experience of mansplaining, including gender-related microaggressions or discriminatory behaviour. The *Well, actually...* corpus is available at https://huggingface.co/datasets/carlaperez/Well_actually_mansplaining.

3.1. Data Collection and Curation

The data has been extracted from a Reddit thread¹ dedicated to telling mansplaining situations experienced by women. The original thread contained 2,914 posts. The authors of this paper have manually pre-processed all the posts in order to select only those stories that expressed mansplaining situations. This step was mainly focused on removing answers or comments to these stories, and parallel conversations or stories that would not be considered to be mansplaining or gender-related microaggressions. We also removed those comments talking about the definition of the concept of mansplaining or personal opinions about it. After this process, the corpus was reduced to 886 paragraphs talking about mansplaining situations.

As an example, the following story of mansplaining is included in the corpus:

"My best friend in college drove a diesel car. EVERY. SINGLE. TIME. we were at the gas station men would either stop what they were doing to shout at us or come running over in a panic to let us know that the green pump is diesel only! [...]"

3.2. Corpus Analysis

The *Well, actually...* corpus is composed of 886 paragraphs written in English. The average length is 96 words, although some of them might contain more than 500 words². The authors of the posts identify themselves mainly as females, although in some cases, the author of the paragraph is self-identified as a male, who is however explaining a mansplaining situation which happened to a woman. In the following subsections, we analyze the corpus in terms of which topics are present (Section 3.2.1), from a more nuanced perspective based on the linguistic and social media features present

¹"Tell me your mansplaining stories, so I can laugh"

²With a standard deviation of 69.55 words.

Topic	Stories
1 Author experience and frustration	109
2 Gender roles and stereotypes	71
3 Mansplaining in tech	58
4 Work experience and Education	56
5 Food use and preparation	43
6 Gender bias at the work place	42
7 Arts, literature and movies	38
8 Pregnancy, childbirth & breastfeeding	33
9 Cars and motor	30
10 Sexual relationships	29
11 Sexism, misogyny	28
12 Health care (physical and mental)	26
13 Fitness and sports	25
14 Gender dynamics	21
15 Gender bias and stereotypes	20
16 Customer service	19
17 Meetings and mansplaining	18
18 Home improv., construction, DIY	16
19 Bra sizes	16
20 Menstruation	15
21 Workplace dynamics & communication	14
22 Languages, grammar, lang. learning	13
23 Coffee	12
24 Accounting, finances, housing market	12
25 Misc. (other topics)	122

Table 1: Main topics and situations where mansplaining occurs in the corpus, ranked by order of prevalence.

in the stories (Section 3.2.2), and for understanding the presence of condescending or misogyny attitudes in the posts (Section 3.2.3).

3.2.1. Topic Analysis

In order to understand better the concept of mansplaining, as well as the situations which are more prone to host this type of discriminatory behaviour, we analyze the main topics covered in the corpus. We use the model Llama-2-70b (Touvron et al., 2023) to extract the main topic for each paragraph. Specifically, we ask the model for the five main topics, as we will explain in Section 4.2, and select its first choice. We then clusterize the topics with BERTopic³ (Grootendorst, 2022), an LM-based topic modeling technique, and assign each paragraph to its corresponding cluster. Table 1 presents the most common topics or situations where we can encounter mansplaining in our corpus, ranked in order of prevalence.

While some of the topics detected by the model refer to the attitudes and feelings of the author (e.g., Topic 1) or gender stereotypes and mansplaining (e.g., Topics 2, 11, 14, 15), in most cases we can observe typical situations where a woman can experience this attitude. It is not unexpected to find

³<https://maartengr.github.io/BERTopic/>

stereotypical men-related areas, such as tech, cars or construction as some of the typical realms of mansplaining. However, we would like to highlight its presence in situations where, inherently, women are supposed to be more experienced, such as pregnancy, childbirth and breastfeeding (Topic 8), bra sizes (Topic 19) or menstruation (Topic 20).

3.2.2. Linguistic and Social Media Features

With the objective of better understanding the concept of mansplaining, we analyze what other biases might be concealed in a mansplaining situation. To this end, we select classification models specialized in social media for analysing features and biases that could be present in some attitudes of the mansplaining situations. In particular, we analyze the mansplaining corpus by applying the following tasks and models:

- 1. Offensive Language Detection.** The fact that a situation is considered mansplaining makes it inherently offensive for the women involved in such situation (Dular, 2021). We analyze to what extent a model trained on detecting offensive language online might identify these posts as being offensive. To this end, we use a model fine-tuned on the SemEval 2019 task on offensive language identification in social media (Zampieri et al., 2019).
- 2. Sentiment Analysis.** We are also interested in analyzing which are the sentiments expressed in a post explaining a mansplaining situation. In particular, we use a simplified sentiment analysis task to classify posts as positive, neutral or negative. The sentiment analysis model is trained on the SemEval-2017 task on Sentiment Analysis in Twitter (Rosenthal et al., 2017).
- 3. Irony Detection.** Given that many of the posts in the corpus are written in a humorous way, we analyze the results of an irony detection model fine-tuned on the SemEval-2018 Irony Detection task (Van Hee et al., 2018). In addition to analyze how many stories are considered to contain irony, we aim at assessing the correlation of the use of this writing style (i.e., the use of irony) with the results provided by other models, such as the offensiveness or the sentiment of those posts.

For all the above tasks we rely on the RoBERTa-based language models trained on Twitter included in the TweetNLP library⁴, which provide competitive and state-of-the-art results in each of the individual tasks (Camacho-Collados et al., 2022). We simply

⁴<https://github.com/cardiffnlp/tweetnlp>

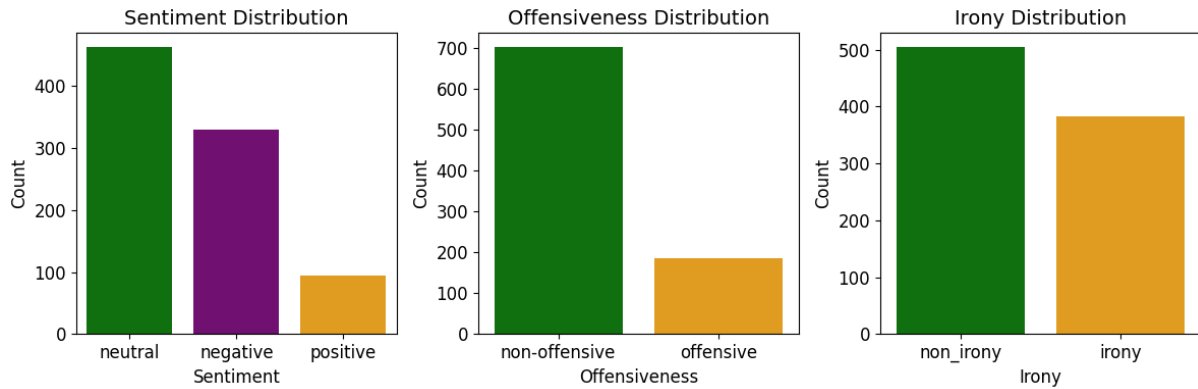


Figure 1: Analysis of sentiment, offensiveness and irony in the positive cases of mansplaining.

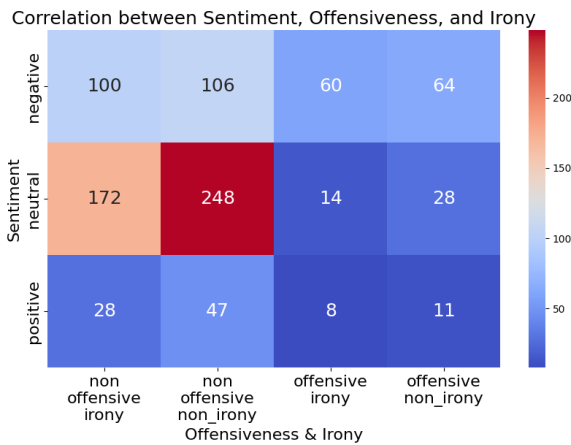


Figure 2: Correlation between sentiment, offensiveness and irony.

use these models to predict their corresponding tasks' labels on each post from the corpus.

Results Figure 1 shows how the models depict the mansplaining stories in the corpus as mainly non-offensive and with neutral sentiment, although a negative sentiment is also common among the mansplaining stories. Both the positive and neutral sentiment detected in these posts might be partially explained by a significant number of ironic messages. In fact, if we analyze the correlation of the results of these three models, as shown in Figure 2, we can observe that there are 222 paragraphs expressing neutral or positive sentiment which were also labelled as being ironic. More surprisingly, there are 295 paragraphs with positive or neutral sentiment which have been labeled as non-offensive and non-ironic. These counter-intuitive results (i.e., mansplaining is offensive in nature and people talking about mansplaining usually do not express a positive sentiment) point to the fact that mansplaining uses a more subtle and indirect type of language and that mainstream models might struggle to understand it.

3.2.3. Condescending Language and Misogyny

In addition to the analysis presented in the previous section, we inspect models with features close to what we could expect from mansplaining posts. In particular, we analyze whether models can identify condescending language or misogyny in the data:

1. **Patronizing and Condescending Language (PCL).** We analyze how a model trained on classifying PCL messages towards vulnerable communities (Pérez-Almendros et al., 2020) might be able to detect the condescending treatment inherent to mansplaining behaviours. This model has been trained on annotated media stories about underrepresented groups.
2. **Misogyny detection.** We use the model released by Attanasio et al. (2022) to assess how many of the mansplaining stories also contain misogyny. This model has been fine-tuned on annotated tweets on misogyny behaviours and attitudes (Fersini et al., 2018).

Results The models that tried to find condescension and misogyny in mansplaining posts throw humble results (none of the posts in the corpus were classified as contained PCL, and only 3.5% -31 posts- were classified as misogynous). On the one hand, the PCL model is not able to find condescension in any of the analyzed posts, despite mansplaining situations containing condescension towards women by definition. This can be due to several reasons such as the type of condescension or the source of the data in which the model has been trained (i.e., news corpora) being significantly different for the model to capture the nuances. On the other hand, we understand that mansplaining is a subtle and sometimes unconscious type of discrimination, so we expected misogyny to be present in the more flagrant cases only. This may include cases where the man was purposely trying to harm

Post	Features
1. I was once invited to the Indian High Commission in London for a drinks event with a client. I invited my husband along as it's probably the fanciest thing I've ever done or will do. Well, the client arrived with his elderly father who, in front of the female, highly poised and clearly Indian high commissioner, proceeded to address my white, under-30 yo, wearing knock off clothes husband as if he was the high commissioner!!! We were appalled but you have to laugh at how some men will do insane mental gymnastics to get things to align to their world view.	Non Offensive, Non Irony, Positive Sent.
2. I was training a new employee (male) and part of the training was sitting in with me when I met with clients. First day, first meeting, I explained to the client that I was training new employee and new employee would be sitting in but not participating in conversation, just listening and learning. Within 30 seconds the client was directing all questions to new employee and ignoring my answers. [...] I asked client if he had a problem with me and would prefer to talk to my boss? Client mansplained to me how males are better at my job because the male brain is better at business because they think without emotion. He asked for the new employee and I told him new employee has been on the job for 2 hours [...]. Client finally accepted I would be the one helping him and appointment continued. After he said - wow you really know your stuff. Of course I know my stuff! I've been in my industry for over 10 years!	Non Offensive, Non Irony, Neutral Sent.
3. Silly woman, there's your problem! Barbeque, often referred to as BBQ, is for MEN. Now stop your silly little plan and get into the kitchen to do the kind of cooking women understand, fire is far too complicated.	Misogyny

Table 2: Selection of examples and their classification by bias detection models.

the woman or women experiencing mansplaining, which are rare cases in the *Well, actually...* corpus.

Table 2 shows some examples of posts for which their classification by the models was unexpected, given the nature of mansplaining attitudes. For instance, the author of example number 2 is clearly offended by the situation, although the model does not capture offensiveness in the post. Furthermore, the positive and neutral sentiment labels of examples 1 and 2 do not correspond to the sentiments evoked by the stories told in them. In the last example, classified as misogyny, the model seems to be missing the irony behind the text. Overall, the analyzed models do not capture many of the subtle nuances contained in mansplaining attitudes, which highlights the necessity to study these types of microaggressions as a specific topic in the NLP community.

4. LLM Evaluation

In the second part of our work, we analyze to what extent LLMs might be able to understand and detect mansplaining situations. We hypothesize that the models might struggle to identify these attitudes due to the need for commonsense reasoning and knowledge of the context that are often required to assess a situation as a case of mansplaining (Lutzky and Lawson, 2019; Koc-Michalska et al., 2021).

4.1. Models

We experiment with three LLMs, two open-source models released by Meta and a third one released

and commercialized by OpenAI. Specifically, we use the following models and variants:

LLaMA-2 (Touvron et al., 2023) are a collection of pre-trained LLMs ranging from 7 billion to 70 billion parameters. In this work, we use the LLaMA models fine-tuned for chat use cases, specifically the models pre-trained with 13 billion and 70 billion parameters.

ChatGPT is a conversational agent powered by GPT-based LLMs. These models are pre-trained on general data extracted from the Internet and fine-tuned to improve results in a variety of downstream NLP tasks, such as question-answering, machine translation or content generation, among many others. Specifically, we use ChatGPT-3.5-Turbo.

4.2. Test Questions

In order to learn more about how LLMs understand mansplaining situations, we select four questions to use as part of a prompt to obtain the models' responses. We experiment with zero-shot learning, as we do not provide a definition of either mansplaining or gender bias, or any instruction to classify the posts. This is because our objective is not to obtain the best possible results for this task, but rather to analyze the understanding and behaviour of models in this matter under their default setting. Therefore, in our research design we decide to opt for zero-shot learning so we do not introduce any biases to the models.

First, we select two questions whose answers can be easily analyzed in a quantitative way.

Question	Model	Yes	No	NA	Acc(%)
Q1: Can you identify gender bias in this situation?	LLaMA-2-13b	696	97	93	78.55
	LLaMA-2-70b	840	10	36	94.80
	ChatGPT	775	110	1	87.47
Q2: Is this a situation of mansplaining?	LLaMA-2-13b	607	231	48	68.51
	LLaMA-2-70b	851	24	11	96.05
	ChatGPT	814	47	25	91.90

Table 3: Answers by LLaMA-2-13b, LLaMA-2-70b and ChatGPT to the questions "Is there gender bias in this situation?" and "Is this a mansplaining situation?" for each one of the posts of the mansplaining corpus. The best results for each question are highlighted in bold.

Specifically, we include the following questions as prompts to the LLMs:

Q1- *'Can you identify gender bias in this situation? Please, answer first with a word: yes or no. Then, can you explain why?'*

Q2- *'Is this a case of mansplaining? Please, answer first with a word: yes or no. Then, elaborate why.'*

Works such as Lutzky and Lawson (2019); Eschmann et al. (2020); Dular (2021) and Koc-Michalska et al. (2021) have related the concepts of gender bias and mansplaining. We ask for both the presence of mansplaining and gender bias, as we hypothesize that the former might be more difficult to identify, while the model might be more familiar with other discriminatory behaviours (Ortiz, 2022; Perez-Almendros and Schockaert, 2022; Naous et al., 2023; Kiehne et al., 2022; Ramezani and Xu, 2023).

Furthermore, we experiment with two open questions which are intended to test whether the models would detect gender-related microaggressions without mentioning them. These questions have the objective of avoiding the potential bias that we can introduce with our prompts to the model (Alnegheimish et al., 2022; Lin and Ng, 2023). The specific questions that we include in the prompt are the following:

Q3- *'What can you infer from this situation?'*

Q4- *'Can you identify up to five topics that appear in the text? Please rank them in terms of how prominent they are in the text'.*

For our experiments, we concatenate each question to each one of the posts contained in the corpus and feed this information to the LLM.

4.2.1. Human Annotation

For Q3 and Q4, the objective is to ascertain to what extent would the models identify mansplaining or other types of gender discrimination without being asked about these phenomena. To evaluate the models' responses, we recruited an annotator to label the models' answers with the following instructions:

For Q3, *'What can you infer from this situation?'*, the annotator was asked to decide whether the LLMs have inferred mansplaining or gender bias in each situation based on their answers. The answer does not need to explicitly mention these specific concepts. For instance, the answer may express that there is an unfair situation related to gender bias towards women. The guidelines provided to the annotator for each of the labels are the following: **0** - *The model does not infer the presence of mansplaining or gender discrimination;* **1** - *Unsure or borderline case;* **2** - *The model infers that there is mansplaining or gender discrimination in that situation.*

For Q4 *'Can you identify up to five topics that appear in the text? Please rank them in terms of how prominent they are in the text'*, we need to annotate if the topic of mansplaining or gender bias/discrimination is present in one of the main 5 topics indicated by the model. The labels for the annotation of this question are as follows: **0** - *The topic is not present in the main five topics;* **1-5** - *The topic is in position 1, 2, 3, 4 or 5;* **6** - *Unsure whether one of the given topics refers to gender discrimination (e.g, the model provides "unfair treatment at the workplace" as a topic, but it does not mention gender).* In the case of the model's response giving more than one topic about gender discrimination or mansplaining, the annotator would select the highest ranked (e.g., if both topics 2 and 4 are about gender bias, they would annotate just "2")⁵.

⁵The annotation guidelines are published at https://huggingface.co/datasets/carlaperez/Well_actually_mansplaining

4.3. Results

In this section, we present the results for both quantitative and qualitative questions of our LLM evaluation.

4.3.1. Quantitative Questions (Q1 and Q2)

Table 3 summarizes the answers of the LLMs when presented with Q1 and Q2, i.e., when asked directly if they can identify either gender bias or mansplaining in each situation. Although the type of language used in these stories is sometimes subtle and some inferences about the context need to be made by the models (e.g., the fact that the person affected by the mansplaining situation is a woman), all three models perform reasonably well when detecting whether there is a discriminatory behaviour towards women. In terms of comparison, LLaMA-2-70b performs better than the other models, correctly identifying mansplaining and gender bias in 96.05% and 94.80% of the stories, respectively. LLaMA-2-13b reports the worst results, which can be explained by the significant smaller size of this model. It also provides more irrelevant answers (NA, in the table), where the model either refuses to answer, gives a blank response or offers a comment which is not related to the question asked. Interestingly, the two largest models (i.e. LLaMA-2-70b and ChatGPT) perform better when identifying mansplaining than gender bias. These results are somewhat counter-intuitive as gender bias is a more general concept. Therefore, our hypothesis was that Q1 would be easier to answer by the models. It is important to emphasize that this evaluation is based on assessing mansplaining situations only, as models may cautiously assess most situations as mansplaining irrespective of whether the actual situation is mansplaining or not.

4.3.2. Qualitative Questions (Q3 and Q4)

The high performance by LLMs in the quantitative answers might be explained by the potential bias introduced in the questions by explicitly asking about the presence of mansplaining or gender bias. In fact, if the models are not asked about these concepts and are given more freedom to answer, as we do with Q3 and Q4, the results are significantly different. For the qualitative questions that required manual annotation (i.e. Q3 and Q4), we only use LLaMA-2-70b and ChatGPT, as these were the best-performing models in the previous experiment. Table 4 shows the results of the models when asked Q3 - "What can you infer in this situation", concatenated to each one of the posts of the corpus. As explained in Section 4.2.1, the human annotator reviewed the models' answers to annotate whether they mentioned something about mansplaining or

Mansplaining or gender discrimination inferred?	ChatGPT (%)	LLaMA-2 70b (%)
No	63.79	53.86
Yes	30.51	40.07
Unclear	5.70	6.07

Table 4: Summary of the results by ChatGPT and LLaMA-2-70b after human annotation. The models were asked "What can you infer from this situation?" in the prompt. This question was concatenated to each one of the posts of the corpus. The annotator then decided whether the models mentioned mansplaining or any type of gender discrimination in their responses.

gender discrimination in their responses. As can be observed, both models often fail to identify any discriminatory behaviour towards women in most cases. However, LLaMA-2-70 detects more cases than ChatGPT, making it a better model to identify subtle gender biases in our corpus.

Table 5 shows the results of the models after asking them Q4 - 'Can you identify up to five topics that appear in the text? Please rank them in terms of how prominent they are in the text', and after annotation. The annotator labelled the answers to assess whether they mentioned the topics of mansplaining or gender discrimination and, if so, in which position. The results show that both models fail to identify these topics as being present at all in many of the situations presented to them. In fact, in the case of ChatGPT, the model considers that most of the stories presented (55.7% of them) do not have any kind of gender discrimination as a prominent topic. Between the two models, LLaMA-2-70b offers a better performance, identifying the problem as a prominent topic (position 1 or 2 in the ranking) in almost 39% of the instances, whereas ChatGPT only identifies gender discrimination as a prominent topic in 28% of the instances.

Qualitative analysis For illustrative purposes, Table 6 includes a couple of examples of posts that were classified as containing both mansplaining and gender bias by the largest models (i.e. ChatGPT and LLaMA-2-70b). However, when asked a less explicit question (i.e., "what can you infer?"), the models are less vehement in their responses. For instance, in the first example, both models highlight and praise the supportive attitude of the woman's boss towards her, but do not seem to see a problem with the discrimination she suffers from the client, as they do not (ChatGPT) or hardly mention it (LLaMA-2-70b). LLaMA-2-70b does allude to *the potential for gender bias* at the end of its explanation, but just to again praise the boss' attitude.

How prominent is mansplaining or gender discrimination in the text?	ChatGPT (%)	LLaMA-2-70b (%)
Not present	55.70	45.77
1	22.61	30.33
2	5.51	8.64
3	5.33	6.25
4	4.41	2.39
5	3.31	2.02
Unclear	3.13	4.60

Table 5: Summary of the results by ChatGPT and LLaMA-2-70b after human annotation. The models were asked "Can you identify up to five topics that appear in the text? Please rank them in terms of how prominent they are in the text", concatenated to each post of the corpus. The annotator annotated whether the models mentioned mansplaining or any type of gender discrimination in their responses and if so, in which position.

In the second example, none of the models identifies mansplaining or gender bias, even though it is explicitly stated in the post. LLaMA-2-70b highlights that "the passenger was trying to make the driver feel safe", and it also mentions that he did not know about the driver's military experience, although the text explains that she shares this background with him. Both models assert that the driver was not interested in the passenger's offer, so she politely dismisses it. These examples illustrate how the models still reproduce and perpetuate gender biases, which can be due to their lack of subtle understanding or commonsense knowledge which is often required to understand such situations.

5. Conclusions

In this paper, we introduce and analyze *Well, actually...*, a corpus of Mansplaining Stories with two main objectives: 1) Improving our understanding of the language used when explaining mansplaining (and other gender related microaggressions), and 2) exploring the capabilities of LLMs to identify and understand this phenomenon. On the one hand, our findings point in the direction of mansplaining using a more subtle and indirect language than other types of biased or discriminatory discourse (such as offensive language, misogyny or condescension). This seems to mislead models devoted to detect these types of biases to misclassify many instances of mansplaining and points towards the need for a model fine-tuned specifically in mansplaining situations to better identify this phenomenon. On the other hand, our experiments with LLMs show that, although they are able to understand and identify mansplaining to some

extent, there is still much room for improvement. For instance, they perform better when the question/prompt is directed towards the answer we want to obtain (i.e., asking them whether the specific discriminatory attitude is present in a given situation). This also highlights the importance of the prompt and how it can affect towards hiding inherent model biases. In fact, when the prompt does not contain gender-related keywords (e.g., when asking the model *what can you infer in this situation* or *which are the main topics*), not only LLMs fail to identify mansplaining, but they also reproduce some gender-related discriminatory attitudes. For instance, they can praise a man who gives unsolicited advice or offers help to women in these situations.

6. Limitations

One of the goals of this paper is to provide the community with a new language resource to motivate further research in mansplaining, still an underexplored research topic in NLP. We have also aimed at better understanding the perception of LLMs regarding this and similar phenomena. However, we are aware that our work presents lacks and limitations that we are confident to solve with (our own and others') future work in this matter.

For instance, our corpus only presents a limited number of mansplaining stories (i.e., 886), all coming from a unique Reddit thread⁶. Although in this paper we focus on the analysis of mansplaining, we believe that further research in mansplaining would benefit from more and more varied data. For instance, fine-tuning a LM for text classification or trying to mitigate mansplaining biases in LLMs would obtain better results with a larger and richer corpus. Nevertheless, the *Well, actually...* corpus offers a starting point to address this challenge.

Moreover, in this work the models' responses were reviewed by only one annotator. At the beginning of the annotation process, we believed that the interpretation of the responses would be easy and not open to much discussion. However, a few instances turned out to not having clear-cut annotations. Nevertheless, the limited resources available (35 paid hours of annotation) made us choose a higher number of annotations over a higher number of annotators. In order to mitigate this limitation, we asked a second annotator to annotate a sample of 5% of the instances. The Inter-Annotator Agreement for this sample is 85%. Both annotators self-identify as European women, with ages between 20 and 40, and with higher education degrees. Please also note that the models' responses have been annotated for the analysis in this paper.⁷

⁶Tell me your mansplaining stories, so I can laugh

⁷The models' responses are available along the

POST: I've been in a call with a company client where we did introductions, I'm the laboratory analyst and was the one working with the stuff they were sending to us. We ran through normal procedure, and I suggested some ways to move forward with their stuff to get the results they wanted. "That's great but what does the man working on our stuff think? (Directed at my boss)" Boss: "the man not working on your stuff thinks the woman working on your stuff is right." Oh the client wasn't happy, but it felt wonderful and big props to my boss for being a legend.

ChatGPT: What can you infer from this situation? From this situation, you can infer the following:

1. The client had a preference for the opinion of the person directly working on their project rather than someone who was not directly involved. 2. The client may have had doubts or concerns about the suggestions made by the laboratory analyst. 3. The laboratory analyst, who is a woman, received positive feedback and validation from their boss, indicating that their suggestions were valued and respected. 4. The boss demonstrated support for the laboratory analyst by standing up for their expertise and opinion, even though they were not directly involved in the project. 5. The client's dissatisfaction suggests that they may have expected a different response or outcome from the call.

Lama-2-70b: What can you infer from this situation? 1. The client had a preconceived notion that a man would be working on their stuff, and was surprised when a woman was working on it instead. 2. The boss was supportive of the laboratory analyst and her expertise, and was not afraid to challenge the client's biases. 3. The laboratory analyst was confident in her abilities and was able to suggest ways to move forward with the client's project. 4. The client was not pleased with the response from the boss, indicating that they may have been expecting a different reaction. 5. The situation highlights the potential for gender bias in the workplace and the importance of having supportive and inclusive leaders who recognize and value the contributions of all team members, regardless of gender.

POST: When I was driving Uber/Lyft for shits and giggles after the military, had a guy say to me "I'll protect you". In reference to a bad neighborhood. This was after I told him I was ex military. He says "I know how to shoot a gun" This guy couldn't whip his way out of a wet paper sack. I informed him that I was proficient with several different weapons and knew self defense on top of being a force protection trainer. Still tried to mansplain, I told him to take a seat.

ChatGPT: What can you infer from this situation? From this situation, it can be inferred that the person driving Uber/Lyft had a conversation with a passenger who made a comment about protecting them in a bad neighborhood. The driver mentioned being ex-military, which prompted the passenger to assert his knowledge of shooting a gun. The driver then asserted their own proficiency with weapons and self-defense, and dismissed the passenger's attempts to explain or assert dominance.

Lama-2-70b: What can you infer from this situation? 1. The passenger was trying to make the driver feel safe by offering to protect them. 2. The driver was not impressed with the passenger's offer and found it amusing. 3. The driver had military experience and was confident in their ability to protect themselves. 4. The passenger was not aware of the driver's military background and was trying to show off their own skills. 5. The driver was not interested in the passenger's offer and politely declined.

Table 6: Examples of posts and models' responses when asked **Q3 - What can you infer from this situation?** Note that these posts were classified as containing gender bias and mansplaining by both models, when asked Q1 - Can you identify gender bias in this situation? and Q2. - Is this a case of mansplaining?

7. Ethical Statement

Mansplaining is an offensive attitude towards women. We are aware that with the release of the *Well, actually...* corpus, as well as with the illustrative examples included in this paper, we are giving a voice to this behaviour. However, our objective is to raise awareness about this and other types of gender-related microaggressions that, unfortunately, are still widely embedded in and accepted by our societies. Moreover, we aim at pointing at the

Well, actually... corpus at https://huggingface.co/datasets/carlaperez/Well_actually_mansplaining

struggles of current NLP systems to detect these discriminatory attitudes and highlight the necessity of 1) being aware, both as researchers and users of AI systems, that these models still reproduce gender biased attitudes, and 2) addressing the necessity of mitigating these biases in NLP systems, in order to obtain more responsible models.

In addition, while this paper is focused on the study of the practice of mansplaining, we do not deny the existence of other types of condescending attitudes, that are not necessarily gender related and that have already been object of other research works in the NLP community.

8. Acknowledgements

Jose Camacho-Collados is supported by a UKRI Future Leaders Fellowship.

9. Bibliographical References

- Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 298–306.
- Sarah Alnegheimish, Alicia Guo, and Yi Sun. 2022. Using natural sentence prompts for understanding biases in language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2824–2830.
- Giuseppe Attanasio, Debora Nozza, Dirk Hovy, and Elena Baralis. 2022. [Entropy-based attention regularization frees unintended bias mitigation from lists](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1105–1119, Dublin, Ireland. Association for Computational Linguistics.
- Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. Tweeteval: Unified benchmark and comparative evaluation for tweet classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650.
- Jeremy Barnes, Laura Ana Maria Oberländer, Enrica Troiano, Andrey Kutuzov, Jan Buchmann, Rodrigo Agerri, Lilja Øvrelid, and Erik Velldal. 2022. Semeval-2022 task 10: Structured sentiment analysis. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, Seattle. Association for Computational Linguistics.
- Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. Semeval-2019 task 5: Multilingual detection of hate speech against immigrants and women in twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.
- Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 2019. [Finding microaggressions in the wild: A case for locating elusive phenomena in social media posts](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1664–1674, Hong Kong, China. Association for Computational Linguistics.
- Caitlin Q Briggs, Danielle M Gardner, and Ann Marie Ryan. 2023. Competence-questioning communication and gender: Exploring mansplaining, ignoring, and interruption behaviors. *Journal of Business and Psychology*, pages 1–29.
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Jose Camacho-Collados, Kiamehr Rezaee, Tayeh Riahi, Asahi Ushio, Daniel Loureiro, Dimosthenis Antypas, Joanne Boisson, Luis Espinosa-Anke, Fangyu Liu, Eugenio Martínez-Cámara, et al. 2022. TweetNLP: Cutting-Edge Natural Language Processing for Social Media. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Abu Dhabi, U.A.E. Association for Computational Linguistics.
- Nicole Dular. 2021. Mansplaining as epistemic injustice. *Feminist Philosophy Quarterly*, 7(1).
- Rob Eschmann, Jacob Groshek, Rachel Chandlerdatt, Khea Chang, and Maysa Whyte. 2020. Making a microaggression: Using big data and qualitative analysis to map the reproduction and disruption of microaggressions through social media. *Social Media+ Society*, 6(4):2056305120975716.
- Emilio Ferrara. 2023. Should chatgpt be biased? challenges and risks of bias in large language models. *arXiv preprint arXiv:2304.03738*.
- Elisabetta Fersini, Francesca Gasparini, Giulia Rizzi, Aurora Saibene, Berta Chulvi, Paolo Rosso, Alyssa Lees, and Jeffrey Sorensen. 2022. Semeval-2022 task 5: Multimedia automatic misogyny identification. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 533–549.

- Elisabetta Fersini, Debora Nozza, and Paolo Rosso. 2018. Overview of the evalita 2018 task on automatic misogyny identification (ami). In *CEUR workshop proceedings*, volume 2263, pages 1–9. CEUR-WS.
- Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2023. [Bias and fairness in large language models: A survey](#).
- Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*.
- Melissa Hall, Laurens van der Maaten, Laura Gustafson, and Aaron Adcock. 2022. A systematic study of bias amplification. *arXiv preprint arXiv:2201.11706*.
- Dirk Hovy and Shrimai Prabhumoye. 2021. Five sources of bias in natural language processing. *Language and Linguistics Compass*, 15(8):e12432.
- Niklas Kiehne, Hermann Kroll, and Wolf-Tilo Balke. 2022. [Contextualizing language models for norms diverging from social majority](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4620–4633, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Karolina Koc-Michalska, Anya Schiffrin, Anamaria Lopez, Shelley Boulianne, and Bruce Bimber. 2021. From online political posting to mansplaining: The gender gap and social media in political discussion. *Social Science Computer Review*, 39(2):197–210.
- Hadas Kotek, Rikker Dockum, and David Q Sun. 2023. Gender bias and stereotypes in large language models. *arXiv preprint arXiv:2308.14921*.
- Anne Lauscher and Goran Glavaš. 2019. Are we consistently biased? multidimensional analysis of biases in distributional word vectors. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 85–91.
- Ruixi Lin and Hwee Tou Ng. 2023. Mind the biases: Quantifying cognitive biases in language model prompting. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5269–5281.
- Ursula Lutzky and Robert Lawson. 2019. Gender politics and discourses of# mansplaining,# manspreading, and# manterruption on twitter. *Social Media+ Society*, 5(3):2056305119861807.
- Thomas Manzini, Yao Chong Lim, Yulia Tsvetkov, and Alan W Black. 2019. Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. In *Proceedings of NAACL-HLT*, pages 615–621.
- Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In *Proceedings of the 12th international workshop on semantic evaluation*, pages 1–17.
- Tarek Naous, Michael J Ryan, and Wei Xu. 2023. Having beer after prayer? measuring cultural bias in large language models. *arXiv preprint arXiv:2305.14456*.
- Roberto Navigli, Simone Conia, and Björn Ross. 2023. Biases in large language models: Origins, inventory, and discussion. 15(2).
- Tolúlopè Ògúnremí, Valerio Basile, and Tommaso Caselli. 2022. Leveraging bias in pre-trained word embeddings for unsupervised microaggression detection. *IJCoL. Italian Journal of Computational Linguistics*, 8(8-2).
- Flor Ortiz. 2022. In-domain and cross-domain classification of patronizing and condescending language in social media and news texts: A study in implicitly aggressive language detection and methods.
- Parth Patwa, Gustavo Aguilar, Sudipta Kar, Suraj Pandey, Srinivas Pykl, Björn Gambäck, Tanmoy Chakraborty, Tamar Solorio, and Amitava Das. 2020. Semeval-2020 task 9: Overview of sentiment analysis of code-mixed tweets. In *Proceedings of the fourteenth workshop on semantic evaluation*, pages 774–790.
- Carla Pérez-Almendros, Luis Espinosa Anke, and Steven Schockaert. 2020. Don't patronize me! an annotated dataset with patronizing and condescending language towards vulnerable communities. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5891–5902.
- Carla Perez-Almendros, Luis Espinosa-Anke, and Steven Schockaert. 2022. [SemEval-2022 task 4: Patronizing and condescending language detection](#). In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 298–307, Seattle, United States. Association for Computational Linguistics.
- Carla Perez-Almendros and Steven Schockaert. 2022. Identifying condescending language: A tale of two distinct phenomena? In *Proceedings of the Second Workshop on NLP for Positive Impact (NLP4PI)*, pages 130–141.

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners.
- Aida Ramezani and Yang Xu. 2023. Knowledge of cultural moral norms in large language models. *arXiv preprint arXiv:2306.01857*.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. Semeval-2017 task 4: Sentiment analysis in twitter. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 502–518.
- Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. Gender bias in machine translation. *Transactions of the Association for Computational Linguistics*, 9:845–874.
- Patrick Schramowski, Cigdem Turan, Nico Andersen, Constantin A Rothkopf, and Kristian Kersting. 2022. Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nature Machine Intelligence*, 4(3):258–268.
- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412.
- Ryan Steed, Swetasudha Panda, Ari Kobren, and Michael Wick. 2022. [Upstream Mitigation Is Not All You Need: Testing the Bias Transfer Hypothesis in Pre-Trained Language Models](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3524–3542, Dublin, Ireland. Association for Computational Linguistics.
- Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1630–1640.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Cynthia Van Hee, Els Lefever, and Véronique Hoste. 2018. Semeval-2018 task 3: Irony detection in english tweets. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 39–50.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Angelina Wang and Olga Russakovsky. 2021. Directional bias amplification. In *International Conference on Machine Learning*, pages 10882–10893. PMLR.
- Zijian Wang and Christopher Potts. 2019. Talkdown: A corpus for condescension detection in context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3711–3719.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. [SemEval-2019 task 6: Identifying and categorizing offensive language in social media \(OffensEval\)](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 75–86, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Hamdy Mubarak, Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin. 2020. Semeval-2020 task 12: Multilingual offensive language identification in social media (offenseval 2020). *arXiv preprint arXiv:2006.07235*.
- Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018. Learning gender-neutral word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Yi Zhou, Jose Camacho-Collados, and Danushka Bollegala. 2023. [A predictive factor analysis of social biases and task-performance in pretrained masked language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11082–11100, Singapore. Association for Computational Linguistics.