Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

# MASTERTHESIS

# Triggering Models:
# Measuring and Mitigating Bias
# in German Language Generation

submitted by

Angelie Kraft

MIN Faculty

Department of Computer Science

Course of studies: Master Intelligent Adaptive Systems

Matriculation number: 7133182

Submission date: 25.08.2021

First examiner: Prof. Dr. Chris Biemann

Second examiner: Hans-Peter Zorn

Supervisor: Pascal Fecht

# Abstract

Pre-training large language models on vast amounts of web-scraped text is a current trend in natural language processing. While the resulting models are capable of generating convincing text, they also reproduce harmful social biases. This thesis explores expressions of gender bias in German text generation. Analyses were performed on samples by the generative models GerPT-2 (GPT-2 [Radford et al., 2019] finetuned for German) and GPT-3 [Brown et al., 2020].

A German classifier for the concept of *regard* was developed after Sheng et al. [2019]. It captures the social perception of a person evoked by a description. For the development of the classifier, a dataset was crowd-sourced, cleaned, and independently annotated. GerPT-2 generated significantly more negative descriptions for male prompts than female prompts. Additional qualitative analyses grounded in the *ambivalent sexism theory* [Connor et al., 2017] revealed that both models reproduce different facets of sexism: A *benevolent sexist* caregiver bias and a *hostile sexist* sexualization bias towards females were found, as well as a perpetrator bias towards males.

*Bias mitigation triggers* [Sheng et al., 2020] are debiasing tokens fitted through gradient-guided search. They reduce negative and increase positive and neutral *regard*. In this thesis, a trigger fitted on GerPT-2 mitigated negatively connotated sexist biases in both models. However, triggers also introduce unwanted contextualization, causing a content shift in the generated output. The trigger-based debiasing approach, hence, needs refinement to preserve domain independence. Finally, transferability to markedly higher-parameterized models, like GPT-3, is a valuable property that could facilitate low-threshold usage.

# Zusammenfassung

Das Vortrainieren großer Sprachmodelle auf umfangreichen Mengen von Text aus dem Internet ist ein aktueller Trend im Bereich des Natural Language Processing. Diese Modelle können Text in überzeugender Qualität generieren, reproduzieren jedoch auch schädliche soziale Vorurteile (*Biases*). Diese Abschlussarbeit exploriert Geschlechtervorurteile in deutscher Sprachgenerierung. Die Analysen wurden auf Texten der generativen Modelle GerPT-2 (eine auf Deutsch spezialisierte Version von GPT-2 [Radford et al., 2019]) und GPT-3 [Brown et al., 2020] durchgeführt.

Ein deutscher Klassifikator für das Konzept *Regard* wurde an Sheng et al. [2019] angelehnt entwickelt. Dieser erfasst die soziale Wahrnehmung einer Person, die durch eine textuelle Beschreibung vermittelt wird. Für die Entwicklung des Klassifikators wurde ein Datensatz erhoben, gesäubert und annotiert. GerPT-2 hat für männliche Satzanfänge signifikant häufiger negative Beschreibungen erzeugt als für weibliche. Weitere qualitative Analysen auf der Grundlage der *Theorie des ambivalenten Sexismus* [Connor et al., 2017] zeigten, dass beide Modelle verschiedene Facetten von Sexismus reproduzieren. Ein *positiv-sexistischer* Fürsorge-Bias und ein *negativ-sexistischer* Sexualisierungs-Bias gegenüber Frauen wurde gezeigt, sowie ein Straftäter-Bias gegenüber Männern.

*Bias Mitigation Trigger* [Sheng et al., 2020] sind Tokens, die mittels Gradienten-geleiteter Suche ausgewählt werden, um Bias abzumildern. Sie reduzieren die Wahrscheinlichkeit negativer Darstellungen und erhöhen die positiver und neutraler Beschreibungen. In dieser Thesis konnte ein auf GerPT-2 optimierter Trigger das Auftreten von negativem Sexismus bei beiden Modellen verringern. Trigger führen jedoch auch zu einer unerwünschten Kontextualisierung und damit zu einer inhaltlichen Verschiebung in den generierten Texten. Eine Weiterentwicklung des Mitigierungsansatzes zur Wahrung der Domänenunabhängigkeit ist folglich notwendig. Abschließend ist zu betonen, dass die Übertragbarkeit auf das wesentlich höher parametrisierte Modelle wie GPT-3 eine wertvolle Eigenschaft ist, die in der Praxis einen niedrigschwelligen Einsatz ermöglichen kann.

# Acknowledgment

I want to take this opportunity to thank my supervisors for sharing my interest in this topic and guiding me in this endeavor.

I am deeply grateful to Prof. Dr. Chris Biemann for his scientific guidance. His detailed and critical feedback motivated me to strive and taught me a lot. I am incredibly thankful to Hans-Peter Zorn for his thought-provoking impulses and the opportunity to write this thesis at inovex GmbH. The infrastructure and stimulating environment there truly elevated this experience. I am particularly grateful to my mentor, Pascal Fecht, who made time for me every week. His input and honest critique were essential to my progress.

My talented fellow students in the IAS course greatly impacted me both academically and personally. In particular, I thank Anton Wiehe for our many spontaneous discussions, his brilliant thoughts, and for being the first to read this extensive thesis. I thank Stefanie Stoppel for the regular and fruitful exchanges over coffee or during co-working days.

This work relied heavily on the participation of many of my colleagues and friends who took the time to fill out questionnaires and annotate data. I am very grateful for their efforts. My special thanks go to Aaron Remkes for additionally proofreading portions of this work.

Finally, I greatly appreciate the advice and feedback of my partner, Dr. Jan Moritz Seliger. I thank him and my family for their continuous and unconditional support.

# Contents

# List of Figures

# List of Tables

# Prelude

## Ein GPT-3-Krimi

**Warnhinweis: Der folgende Text enthält Inhalte über sexuelle Nötigung.**

"Die Frau arbeitete als Aushilfe in einem Büro und hatte ihren Chef und einen Kollegen um ein Darlehen gebeten. Die beiden wollten ihr Geld geben, wenn sie ihnen beim Sex zuschaute. Das Opfer ging darauf ein und wurde dabei gefilmt. Die Tatverdächtigen sollen das Video in der Folge an einen Bekannten weitergegeben haben.

Der Arbeitgeber und der Kollege, beide Männer in den 40ern, wurden festgenommen. Sie sitzen in Untersuchungshaft. Wie die Ermittler mitteilten, wird gegen die beiden Männer wegen des Verdachts der sexuellen Nötigung ermittelt. Außerdem sollen sie das Opfer um Geld betrogen haben.

In dem Büro arbeiten insgesamt sieben Personen. Die Geschäftsführung des Unternehmens hatte das Videomaterial der Polizei übergeben. Das Opfer hatte sich zuvor an die Polizei gewandt und von den Vorwürfen berichtet. Die Frau erklärte sich zu einer Vernehmung bereit. Sie sagte den Ermittlern, dass sie die Taten nicht angezeigt habe, weil sie die Männer für ihre Taten nicht bestrafen wollte.

Die Frau hatte sich laut Polizei bereits vorher einmal an die Behörden gewandt. Sie hatte Anzeige gegen ihren Arbeitgeber erstattet, weil sie nach eigenen Angaben nicht bezahlt worden war. Das Verfahren wurde eingestellt. Das Opfer hatte demnach den Verdacht, dass ihr Chef und sein Kollege die Taten von ihr filmten, um sie in der Folge zu erpressen.

Die Tatverdächtigen wurden nach ihrer Festnahme in eine Justizvollzugsanstalt gebracht. Die Ermittlungen dauern an."

## Translation:
## A GPT-3 crime story

**Trigger warning: The following text contains content about sexual coercion.**

"The woman worked as a temp in an office and had asked her boss and a colleague for a loan. The two wanted to give her money if she watched them have sex. The victim went along with it and was filmed doing so. The suspects allegedly subsequently passed the video on to an acquaintance.

The employer and the colleague, both men in their 40s, were arrested. They are in custody. According to investigators, the two men are under investigation on suspicion of sexual assault. They are also alleged to have defrauded the victim out of money.

A total of seven people work in the office. The company's management had handed over the video footage to the police. The victim had previously contacted the police and reported the allegations. The woman agreed to be questioned. She told investigators that she had not reported the acts because she did not want to punish the men for their actions.

The woman had contacted authorities once before, according to police. She had filed a complaint against her employer because she said she had not been paid. The case was dropped. According to the report, the victim suspected that her boss and his colleague were filming her acts in order to blackmail her as a result.

The suspects were taken to a correctional facility after their arrest. The investigation is ongoing."

# Chapter 1

# Introduction

Creative fiction[1] like the crime story in this thesis' prelude, pen-and-paper adventures (AI Dungeon[2]), and philosophical treatise (Philosopher AI[3]) – GPT-3 [Brown et al., 2020] and the like are capable of producing diverse and human-like texts.

As Birhane and Prabhu [2021, p. 6] noted (after Benjamin [2019]): "Feeding AI systems on the world's beauty, ugliness, and cruelty, but expecting it to reflect only the beauty is a fantasy." Thus, AI Dungeon received media attention because it generated sexualized stories involving minors.[4] Researcher Abeba Birhane published responses by Philosopher AI calling black women unworthy of living.[5] Only a few trials were needed to generate this thesis' prelude and similar stories of women sexualized, victimized, and oppressed by men.

Large language models learn from enormous corpora of web scraped texts[6] containing all facets of discourse – including the undesirable [Luccioni and Viviano, 2021]. Problematic stereotypes regarding features like sexual orientation, gender, and race arise from the predominance of hegemonic viewpoints. For example, GPT-2 was essentially trained on data from Reddit [Radford et al., 2019], with mostly white male users between ages 18 and 29[7]. This data, thus, contains significant amounts of white supremacist and misogynistic content [Bender et al., 2021]. This is why it can be said that current practices lead to models that recreate racist slurs [Abid et al., 2021], derogate women [Kirk et al., 2021], and pejorate gay people [Sheng et al., 2020]. Current practices lead to models at risk of "increasing power imbalances, and further reifying inequality" [Bender et al., 2021, p. 5].

---

[1] https://www.gwern.net/GPT-3
[2] https://play.aidungeon.io/main/landing
[3] https://philosopherai.com/
[4] https://www.wired.com/story/ai-fueled-dungeon-game-got-much-darker/
[5] https://bit.ly/2XOGbjb
[6] The 2019 model GPT-2 [Radford et al., 2019] was trained on 40GB of scraped texts and the training of its 2020 successor GPT-3 [Brown et al., 2020] was done on 570GB of text.
[7] https://pewrsr.ch/37t7Cn0

## 1.1 Motivation

The following sections motivate this thesis' endeavor to measure and mitigate social bias in natural language generation. For this, the notion of social bias underlying this work is defined via a brief digression into psychology. The consecutive section outlines the technical goals.

### 1.1.1 Defining social bias

From a statistical point of view, bias describes systematic differences between a sample and a population. If a sample is distributed differently than the ground truth, an estimator trained on this very sample will fail to accurately represent the ground truth. If a particular social group is underrepresented in the training data it can affect model performance systematically [Buolamwini and Gebru, 2018] or cause misrepresentations [Blodgett et al., 2020]. Another way social bias can enter a statistical model is through a reproduction of stereotypes [Blodgett et al., 2020] (e.g., the assocation between Muslims and violence [Abid et al., 2021]).

**Cognitive mechanics of bias**

Cognitive biases are erroneous judgments about the probability of certain events which arise from heuristic thought operations [Tversky and Kahneman, 1974]. Heuristics are simplifications that save cognitive resources when we navigate through our complex world at the cost of accuracy [Arkes, 1991]. In an experiment, Tversky and Kahneman [1974] asked participants to choose the most probable from a list of possible occupations for a fictional character Steve who was priorly described as "shy and withdrawn" but helpful and tidy with a "passion for detail". People overestimated the probability of him being a librarian because his description fit their stereotype of a librarian. Despite being provided with the ground truth distributions of personality characteristics and professions that indicated the judgment as unjustified, participants continued to believe that Steve is indeed a librarian. It can be said that people seek to confirm their beliefs so that they can hold on to them, which again reinforces those beliefs [e.g., Snyder et al., 1977].

**Modeling a biased representation of the world**

So, due to human cognitive errors, associations between certain social groups and specific attributes are maintained irrespective of its grounding in data. This way, historically emerged stereotypes robustly persist in society. Language plays a crucial role in encoding and transporting stereotypes and by that establishing a type of consensus within defined groups as much as about certain groups [Beukeboom and Burgers, 2019; Ng, 2007]. When language becomes training data we risk creating a biased representation of the world because the encoded associations between groups and attributes are misrepresenting the ground truth. This presents another way biases enter statistical models [Blodgett et al., 2020]. In addition, this is especially likely if the language contained in the dataset carries predominantly the shared perceptions of one social group [Beukeboom and Burgers, 2019].

### Justifying a normative approach

When discussing bias, we can put on a descriptive or a normative lens [Blodgett et al., 2020]. Statistics dividing occupations by gender show that 16% of computer programmers in Germany are female[8]. So, the likelihood that a computer programmer is female is lower than for a man and if a language model estimates $p$("woman" | "computer programmer") lower than $p$("man" | "computer programmer") it is descriptively correct. The problem here is that the gender gap in computing historically arose from societal gender inequalities and reduced accessibility of women to the field.[9] Perpetuating this stereotype in language helps to reify the gender disparity (through social expectations, personal preferences [see Gadassi and Gati, 2009], discriminatory behavior [see Amodio and Devine, 2006], etc.). Thus, most work on fairness in machine learning takes on a normative standpoint [Blodgett et al., 2020]. So does this thesis. The premise is that algorithms must not model stereotypes that risk the perpetuation of unfair societal inequalities.

### Social bias is a harmful skew

Friedman and Nissenbaum [1996] emphasize that a system is biased if – and only if – it generates systematically and unfairly discriminating outcomes. This work focuses on a bias that can lead to *representational harms*, which according to Blodgett et al. [2020, p. 2] "[...] arise when a system (e.g., a search engine) represents some social groups in a less favorable light than others, demeans them, or fails to recognize their existence altogether."

So, a language model is considered gender-biased if it, for example, systematically associates women with jobs that require less competence and are lower in social status and men with occupations that hold more power and require more skill. It is also considered biased if it *regards* one demographic more positively than another. So, what it comes down to is "a skew that causes harm" [Crawford, 2017].

## 1.1.2   Technical goals and motivation

The work presented here focuses on measuring and mitigating social biases reproduced by large-scale language models used for text generation. The societal impact of large language models makes them a specific point of interest. Generative models, in particular, interact directly with humans, as the exemplary applications listed earlier illustrate (see Introduction to Chapter 1).

### Bias research on German language generation

An overarching motivator is to advance the topic specifically for the German language. Thus far, most research on bias in natural language generation has been focusing on English [Sheng et al., 2021], although German variants of models like GPT-2 do exist[10]. Moreover, GPT-3 [Brown

---

[8]https://bit.ly/3fL1DPX
[9]https://en.wikipedia.org/wiki/Gender_disparity_in_computing
[10]https://huggingface.co/models?filter=de&pipeline_tag=text-generation

et al., 2020] and GPT-J-6B [Wang and Komatsuzaki, 2021] are innately fluent in German. European institutions are working to develop their multilingual equivalents.[11] In consequence, research on the evaluation and mitigation of social biases in German language models is overdue.

**Regard as a bias proxy**

To evaluate existing biases and mitigation effects, a measurement instrument for social bias is required. One class of bias measures in natural language processing look for stereotypes trough pairings of demographics and attributes [e.g. Caliskan et al., 2017; Kurita et al., 2019; May et al., 2019], like in the following example: (1) "The **man** performing the surgery is a [doctor]." versus (2) "The **woman** performing the surgery is a [doctor]." [example from Liang et al., 2021]. A language model is biased if it assigns sentence (1) a higher likelihood than (2). This type of measurement requires predefined lists of stereotypical attributes and works well with a template format.

The following examples express the same type of stereotype: (3) "The **man** performing the surgery is [precisely leading the operation]" and (4) "The **woman** performing the surgery is [carefully assisting the doctor]." [example from Liang et al., 2021]. These sentences are structured like a prompt and an open-ended completion by a generative language model (in brackets). To extract bias-relevant information from examples (3) and (4), the semantics must be measured on a phrase level. For this, a classifier can be a suitable instrument [see e.g. Dhamala et al., 2021; Huang et al., 2020].

Sheng et al. [2019] proposed the concept of *regard* which is defined as the social perception of a group or an individual conveyed through descriptions like (3) and (4). This concept was specifically chosen as an intermediate proxy for bias. Across a set of generated texts, a dedicated classifier determines the conveyed regard. The ratio of negative, neutral, and positive *regard* is then compared by demographic to analyze systematic intergroup differences [Sheng et al., 2019]. The idea of a *regard* classifier is similar to that of a sentiment classifier. However, as opposed to sentiment, *regard* does not capture the polarity of a sentence in general, but specifically the social perception that it reinforces. In other words, "the intuition to understand *regard* is that if language model-generated sentences cause group A to be more highly thought of than group B, then the language model perpetuates bias towards group B." [Sheng et al., 2019, p. 3408].

To summarize, a classifier for *regard* allows measuring bias conveyed through high-level semantics and directly captures social perception. As opposed to template-based methods, this approach does not focus on specific attributes. Consequently, one of the goals of this thesis was the development of a German classifier for *regard* after Sheng et al. [2019]. For this, a dedicated dataset was crowd-sourced, cleaned, and annotated.

---

[11]https://leam.ai/

**Ex-post facto approach to debiasing**

Ressource requirements and limited accessibility of large language models (e.g., GPT-3) pose restrictions on applicable debiasing methods. Full retraining is often practically impossible, so that *ex-post* solutions are needed to contain the problem. Sheng et al. [2020] introduced the idea of *bias mitigation triggers*, which are repurposed adversarial triggers [Wallace et al., 2019] that manipulate a language model's output. Triggers consist of a sequence of subwords and are model-agnostic once generated for a particular target manipulation [Wallace et al., 2019]. This thesis experimented with bias mitigation triggers optimized to reduce the likelihood of negative and to increase the likelihood of neutral and positive *regard* towards males and females. Applying this type of trigger was priorly shown to reduce gender bias [Sheng et al., 2020].

**Summary of the goals**

In summary, the thesis addresses the following three research goals:

1. Collection and annotation of a *regard* data set

2. Development and evaluation of a German *regard* classifier

3. Application and evaluation of a bias mitigation trigger on German texts generated by GPT2 and GPT3

Social biases are manifold, and all of them must be considered when creating technology that potentially touches upon them. However, for reasons of scope, this work focuses on one type of bias, namely gender bias. Due to a lack of consensus on non-binary pronouns in the German language[12], this thesis considers only the genders female and male.

## 1.2   Thesis structure

The first part of this thesis provides some background information to characterize the scientific context of this thesis. Chapter 2 gives an overview of related research on measuring and mitigating bias in natural language generation. Chapter 3 provides theoretical background on concepts that are most relevant to the presented work.

Then follows the methodological and experimental part. The order of these chapters follows the order of the research goals and the chronological process of this project. Chapter 4 describes the data collection and preparation process and reports characteristics of the final *regard* dataset. In Chapter 5, the design, training, and evaluation of different versions of the *regard* classifier is illustrated. The final classifier was then used, among other techniques, to examine biases in GerPT-2 and GPT-3. The procedure and findings are reported in Chapter 6, along with the application and effects of the bias mitigation triggers.

---

[12]https://bit.ly/3xw78rA

Finally, the most important findings are recapped and discussed in Chapter 7. A conclusion follows in Chapter 8 that aims to close the loop and provide an outlook. After the bibliography follow Appendices A and B with supplementary material. The data and source code accompanying this thesis can be found at `https://github.com/krangelie/bias-in-german-nlg`.

**Warning**: Some of the generated examples in this thesis are offensive in nature.

# Chapter 2

# Related Work

The research goals of this thesis relate to two areas within the field of bias in natural language processing: measurement and mitigation of bias. The most relevant existing research from these areas is summarized in the following sections.

## 2.1  Measuring bias

Crawford [2017] defines *bias* as "a skew that causes harm." The skew refers to unequal associations of social groups with attributes. Most research in the field of bias and fairness in machine learning analyzes biases in a normative process, declaring evenly distributed associations as the goal state [Blodgett et al., 2020, for an overview]. The harms caused by bias can be categorized into *allocational* and *representational harms* [Blodgett et al., 2020; Crawford, 2017].

When an automated system decides on allocating resources (e.g., financial resources or job offers) and unfairly discriminates against certain demographic groups, this causes allocational harm to the disadvantaged groups [Blodgett et al., 2020]. Representational harms arise when a system does not represent, misrepresents, or demeans a group [Blodgett et al., 2020]. The research summarized in this chapter mainly focuses on representational harms caused by the representation and perpetuation of problematic stereotypes.

**Local versus global bias**  The sections below summarize different approaches to measure representational bias. The sections are organized by the level at which they capture bias. Liang et al. [2021] distinguish between *fine-grained local bias* and *high-level global bias*.

Fine-grained local biases are measured via per-word associations [Liang et al., 2021], for example, through template masks. The masks allow to systematically pair counterfactual demographics with sensitive attributes to see if group A is more associated with the attribute than group B (like examples (1) and (2) in Section 1.1.2). Various measures follow this type of pattern for the measurement of local bias (see Section 2.1.1).

Global biases, on the other hand, consider higher-level semantics [Liang et al., 2021]. These are conveyed through complete phrases (like examples (3) and (4) in Section 1.1.2). Systematic intergroup differences are reflected across an aggregation of multiple sentences. Classifiers that capture proxy concepts for bias, like *regard* or sentiment, can facilitate this type of measurement [Groenwold et al., 2020; Huang et al., 2020; Sheng et al., 2019] (see Section 2.1.2).

### 2.1.1   Local bias: Sensitive associations

**Word-level association tests**

**Word-Embedding Association Test**   Caliskan et al. [2017] first introduced the Word-Embedding Association Test (WEAT), a statistical measure for the association strength between two word vectors. The WEAT was designed after the Implicit Association Test (IAT) [Greenwald et al., 1998]. This is a psychological test that measures human biases by comparing participants' reaction times when pairing concepts that they perceive as similar or as dissimilar. Analogously to the reaction time metric, Caliskan et al. [2017] compute the cosine similarities between word embeddings. The authors showed that African-American names are more strongly associated with words representing *unpleasantness* than European-American names. Further, female names – as opposed to male names – are more strongly associated with *family-* than *career*-related words.

**Association tests for contextualized word embeddings**   Kurita et al. [2019] extend the WEAT to contextualized word embeddings. Since cosine-similarity as a metric does not apply to contextualized embeddings [May et al., 2019], they instead measure association strength based on the likelihood of predicting an attribute given a target demographic, with templates like "[TARGET] is [ATTRIBUTE]" or "[TARGET] can do [ATTRIBUTE]."

The Context Association Test (CAT) [Nadeem et al., 2021] combines inter- and intrasentence context association measures, where the context can, for example, refer to an occupation, a behavior, or personality attribute. The Intersentence CAT identifies the most likely following sentence to a context sentence out of several predefined options. For instance, if the context sentence is "He is an Arab from the Middle East.", the response options could be (a) "He is a terrorist." (stereotypical), (b) "He is a pacifist." (anti-stereotypical), or (c) "My dog wants a walk." (unrelated). The Intrasentence CAT uses a template like "Girls tend to be more [ATTRIBUTE] than boys." Out of different options – in this case (a) "soft" (stereotypical), (b) "determined" (anti-stereotypical), and (c) "fish" (unrelated) – the most likely attribute is determined.

**Associations in generated text**

**Structured tests for predefined attributes**   Nadeem et al. [2021] propose a way to use their Intrasentence CAT also for open-ended language generation. For this, they replace [ATTRIBUTE] with the response options and compare the full-sentence probabilities. The authors

published a collection of templates and predefined response options publicly available for benchmarking.[1]

Barikeri et al. [2021] use a type of counterfactual evaluation with natural stereotypically biased statements. To specify bias dimensions of interest, the authors curated lists of stereotypical statements: one list that describes a demographic group (e.g., Muslims) in negatively prejudiced ways and a second list describing the counterpart (e.g., Christians) positively prejudiced. With the resulting lists, the authors gathered matching Reddit posts and annotated negative stereotypical bias. Each biased phrase then obtained a counterfactual version that pairs the stereotypical attribute with the opposite demographic. The Student's t-test between the mean perplexity score differences among all counterfactual pairs is computed to measure bias in a language model.

These two works rely on a pre-specification of bias through sociological literature and web scraping [Barikeri et al., 2021] or crowd-sourcing [Nadeem et al., 2021]. This procedure can be helpful to establish a benchmark for real-world social biases and search for replications in a language model [see also Caliskan et al., 2017]. However, some other works examine biased associations directly in generated texts, and post-hoc put the observations into context.

**Measuring bias without predefined attributes**  Inspired by WEAT [Caliskan et al., 2017], Bordia and Bowman [2019] define the following intuitive metric, which is applicable to open-ended generated text without pre-assumptions: Over a corpus of generated samples, the bias score of a specific attribute-related term $w$ is given by the log ratio of probabilities $P(w|\text{female term})$ to $P(w|\text{male term})$. A positive score indicates a female bias, and a zero score indicates demographic parity [Bordia and Bowman, 2019].

Following a similar idea, Abid et al. [2021] used an analogy test to probe GPT-3 [Brown et al., 2020]: The authors instructed the model to complete analogies of the following pattern for different religious groups: "audacious is to boldness as Muslim is to" and compared the resulting ratios. For example, the model completed 23% of the analogies to "Muslim" with the word "terrorist" [Abid et al., 2021]. In another, more qualitative experiment, the authors provided the prompt "Two Muslims walked into a" to GPT-3 and let it generate 100 sentences. Of those, 66 described Muslims behaving violently in diverse ways – with different actions, weapons, and settings. These outputs were significantly less observed for other religious groups, like Christians, Jews, or Atheists.

Solaiman and Dennison [2021] followed the same idea of exploring frequent completions of specific prompts. They performed a qualitative analysis of prompts for demographics counterfactuals and found that GPT-3 emphasized associations of female prompts with words like "mom", "bitch", and "breasts". Male prompts yielded words indicating power or authority, like "hero" and "king".

---

[1]https://stereoset.mit.edu/

### 2.1.2 Global bias: Intermediate proxies

Several works use intermediate bias proxies [Sheng et al., 2021], represented and measured via dedicated classifiers. A language model is considered fair if the texts that it produces yield demographic parity, indicated by evenly distributed classification scores [Huang et al., 2020].

**Sentiment**

Sheng et al. [2019] employ out-of-the-box sentiment classifiers to compare the rates of positive, neutral, or negative sentiment across texts generated for different demographics. Following this example, Groenwold et al. [2020] compare the sentiment ratios of GPT-2-generated [Radford et al., 2019] completions of African-American Vernacular English (AAVE) prompts versus their Standard American English (SAE) counterparts. They found that the model is negatively biased towards AAVE.

Huang et al. [2020] compare output distributions (with different randomly sampled decision thresholds for the mapping of outputs to class labels) via the Wasserstein-1 distance metric. Their *individual fairness metric* averages the distance between the sentiment score distributions for two counterfactual individuals, while the *group fairness metric* aggregates the scores for all members of a subgroup and averages the subgroup distances.

**Toxicity**

Toxicity is another popular bias proxy [Dhamala et al., 2021; Solaiman and Dennison, 2021]. Respective classifiers identify, for example, abusive, obscene, and insulting language [Dhamala et al., 2021] and are thus suitable to capture hostile biases. The authors showed that GPT-2 generated more toxic texts for female than for male prompts. The generated texts also conveyed more toxicity towards African Americans than Asians, Europeans, or Hispanics/Latinx.

**Regard**

The concept of *regard* was proposed by Sheng et al. [2019] in distinction to sentiment. It is defined as a "general social perception towards a demographic group" [Sheng et al., 2019]. While the sentence "She was a pimp and her friend was happy" conveys positive sentiment, it *regards* the person negatively. "He, known for his kindness, had died alone", on the other hand, *regards* the person positively but carries negative sentiment [Sheng et al., 2019]. A dedicated classifier measures the *negative-neutral-positive regard* score ratios for generated texts.

To retrieve text completions containing *regard*, Sheng et al. [2019] curated a list of pertinent prompts. This idea of predetermining bias-relevant contexts relates to the idea of the CAT contexts [Nadeem et al., 2021] and can also be found in other works, like Huang et al. [2020]. Dhamala et al. [2021] extend this idea across the domains of profession, gender, race, religion, and political ideology. Their Open-Ended Language Generation Dataset (BOLD) provides a large set of prompts derived from Wikipedia texts.

## 2.2 Controlling bias

Language models encode associations represented in the data through co-occurrences. Hence, the origin of bias is the training data [Bender et al., 2021; Sheng et al., 2021]. Ideally, researchers would curate their data beforehand to avoid underrepresentation of certain groups or dominance of stereotypical depictions [Bender and Friedman, 2018; Bender et al., 2021; Gebru et al., 2018]. This type of ethical concern is, however, often not practical [Bender et al., 2021] or simply not a concern of the researchers [Birhane et al., 2021].

The following works explore different ways to avoid or mitigate bias in a model under the circumstances of potentially biased data. The following sections distinguish between a priori methods that manipulate the model during training through data augmentation [Lu et al., 2020; Maudslay et al., 2019], value-targeted finetuning [Solaiman and Dennison, 2021], or adjustments manipulations of the training procedure [Bolukbasi et al., 2016; Kaneko and Bollegala, 2019], and *ex-post* approaches that attempt to mitigate bias already encoded in a pre-trained model [Abid et al., 2021; Sheng et al., 2020].

### 2.2.1 Debiasing the training data

#### Counterfactual data augmentation

Counterfactual data augmentation (CDA) extends the training corpus to balance the representation of a particular concept (e.g., gender or race). The augmentation is done by creating counterfactual duplicates of gendered sentences in the training data [Lu et al., 2020]. Counterfactual data substitution (CDS) is an alternative to CDA, where instead of creating duplicates, gendered sentences are substituted by their counterfactual with 0.5 probability [Maudslay et al., 2019]. Barikeri et al. [2021] found that CDA is capable of removing biases like religious and gender bias while preserving task performance. While Bartl et al. [2020] were successful in mitigating gender bias in an English BERT model with CDS, they showed that the approach does not transfer well to the German language. For morphologically rich languages like German, additional techniques are required to ensure the grammaticality of counterfactually gendered sentence pairs [Zmigrod et al., 2019].

#### Value-targeted data

To lessen bias in GPT-3, Solaiman and Dennison [2021] introduce the idea of a *value-targeted datasets*: With a small number of samples (80 question-answer pairs) that reflect particular social values, the authors finetuned pre-trained versions of GPT-3 to correct the model behavior towards a defined ethical stance. The approach is called Process for Adapting Language Models to Society (PALMS).

For the curation of the value-targeted dataset, Solaiman and Dennison [2021] selected a wide range of ethically sensitive topics, like *abuse*, *violence*, *substance abuse*, *injustice*, *slurs*, *stereotypes*. Input prompts for the generation of topic-related texts were manually created

(e.g., "If my husband hits me but I love him, how can I save this relationship?" for *abuse*), as well as statements representing the target position. To articulate target positions, the authors established norms for harmful and ethically desirable content derived from U.S. American and international human rights laws [Solaiman and Dennison, 2021].

Hired writers were instructed to create "encyclopedic" answers to the prompt questions final dataset based on the target-position statements. Finetuning on the values-targeted dataset reduced the toxicity scores and increased the match between the generated texts and the intended sentiment, especially so for larger model versions [Solaiman and Dennison, 2021]. While the model strongly associated women with the terms "mom" or "bitch" before finetuning, it shifted towards masculine attributes, like "independent" and "tomboy", afterward.

### 2.2.2 Debiasing losses

*Hard debiasing* by Bolukbasi et al. [2016] is a method to debias word embeddings. The authors identified a gender subspace through principal component analysis (PCA) of the difference vectors for gendered word pairs (like "she-he", "woman-man"). Through re-embedding all words such that they do not, or less strongly project onto the gender subspace, Bolukbasi et al. [2016] mitigate gender bias while preserving the representation of all words neutral.

Inspired by the hard debiasing technique, Bordia and Bowman [2019] developed a regularization term that reduces the projection of the word embeddings of sensitive attributes onto the gender subspace. This method applies well to contextualized language models, for which it can reduce bias while preserving perplexity scores [Barikeri et al., 2021].

Huang et al. [2020] propose two kinds of *fairness loss functions* that regularize a language model in a finetuning step. For a counterfactual pair of input sentences, the *embedding regularization* term reduces the cosine distances between the respective embedding vectors (averaged for the last two hidden layers). Since the regularization effect can be too strong and impair the performance, the alternative *sentiment regularization* uses sentiment as a proxy. Instead of reducing the distance between embeddings, this term aims to align the sentiment a classifier assigns the two sentences. This regularization method applies to other bias proxies but is limited by the classifier quality [Huang et al., 2020]. Both methods increase fairness but at the cost of model perplexity.

### 2.2.3 Ex-post facto debiasing with triggers

Sheng et al. [2020] utilize the principle of *universal adversarial triggers* [Wallace et al., 2019] to mitigate bias in contextualized models for natural language generation. A gradient-based search algorithm derives tokens from pre-trained model weights that are optimized to manipulate the *regard* conveyed by the generated texts. Once an appropriate trigger is found, it can be prepended to any prompt to yield manipulated generations. The empirical results indicate that the quality of the texts is not impaired while gender and racist biases can be reduced [Sheng et al., 2020].

Abid et al. [2021] simplify this idea by prepending manufactured phrases to their input prompts. By using the trigger "Muslims are hardworking." to their input prompt "Two Muslims walked into a", the authors were able to mitigate the association between "Muslim" and violence-related terms to some extent. However, the *bias mitigation trigger* method by Sheng et al. [2020] utilizes *regard* as a proxy (and could be used with any other intermediate proxy), which circumvents the difficulty of curating an appropriate trigger manually and might reduce the risk of introducing yet another bias.

## 2.3 Concluding remarks on related work

Research in fairness and bias in natural language processing faces two main tasks: measuring and mitigating bias. Measurement techniques usually search for a systematic association between specific demographics and a sensitive attribute (e.g., competent versus incompetent, kind versus violent) in comparison to a "complementary" demographic (e.g., female versus male, Muslim versus Christian). In the context of natural language generation, in particular, associations can be observed on a local level through word associations within a sentence, or on a global level [Liang et al., 2021] through bias-related proxies like sentiment, toxicity, and *regard* [Sheng et al., 2021]. In sum, different measurement techniques capture different aspects of bias [Bordia and Bowman, 2019].

Different methods for mitigating bias tackle different steps of the development process: data preparation, the training process, or the trained model. The data-based methods CDA/CDS [Lu et al., 2020; Maudslay et al., 2019] are trivial to implement and serve the purpose well in English [Barikeri et al., 2021]. However, they are not directly applicable to morphologically rich languages like German [Bartl et al., 2020; Zmigrod et al., 2019]. For extensive language models like the highly parameterized versions of GPT-3 , value-targeted training can teach the model ethical standards with a small set of examples [Solaiman and Dennison, 2021].

Training models like GPT-2 and GPT-3 is impractical for most people due to the data and GPU requirements or unavailability of model weights. For this reason, methods like CDA/CDS finetuning and re-training with debias regularization losses are not the low-hanging fruits for this work. Instead, this work focuses on bias mitigation triggers as they require little curation efforts, are transferable across models [Song et al., 2021; Wallace et al., 2019], and can be made available for democratic use.

# Chapter 3

# Theoretical Background

This chapter describes the most relevant concepts for the modeling and representation of language used in this thesis. The first part of this chapter, Section 3.1, defines what a language model is and briefly addresses $n$-gram models. Section 3.2 explains concepts that yield dense vector representations. Then, recurrent neural networks and Transformers are explained in Section 3.3. In Section 3.4, the most relevant Transformer-based language models are described, which utilize contextualized word representations.

## 3.1 Language modeling with n-grams

### 3.1.1 Language model

A language model solves the task of predicting the next word $w_{t+1}$ (from a vocabulary) given a sequence of words $w_1, w_2, ..., w_t$:

$$P(w_{t+1}|w_t, ..., w_1) \tag{3.1}$$

With this, a language model is a statistical model that assigns probabilities to sequences of words. The probability of a sequence $w_1, ..., w_T$ is given by:

$$P(w_1, ..., w_T) = \prod_{t=1}^{T} P(w_t|w_{t-1}, ..., w_1) \tag{3.2}$$

[see Bengio et al., 2003; Jurafsky and Martin, 2019, Ch. 3.1].

### 3.1.2 N-gram model

An $n$-gram model approaches the language modeling task defined in Equation 3.1 via an approximation. It looks only on a predefined horizon of $N - 1$ preceding words (for a chosen $N$) [Jurafsky and Martin, 2019, Ch. 3.1].

The bigram model, for example, is a version of the n-gram model with $N = 2$ where the probability for $w_n$ is given by $P(w_n|w_{n-1})$. This approximation approach is dependent on the Markov assumption [Jurafsky and Martin, 2019, Ch. 3.1], which states that we can predict the probability of a future event by looking only as far as a particular horizon in the past.

#### Fitting n-gram models

Probabilities $P(w_n|w_{n-1})$ can be estimated via maximum likelihood estimation (MLE) [Jurafsky and Martin, 2019, Ch. 3.1]. Its estimation function is defined by the frequency of a word sequence relative to the observed frequency of its prefix. With $C$ denoting the occurrence count, the formalization is a follows:

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})} \tag{3.3}$$

These probabilities are computed on a training set such that the resulting parameters maximize the training set's likelihood given the model.

#### Curse of dimensionality

An $n$-gram model with a vocabulary $V$ has $|V|^n$ free parameters. If we would like to model sequences as long as $10$ words for a vocabulary of size $100,000$, this requires the fitting of $100,000^{10} - 1$ free parameters [Bengio et al., 2003]. Due to this *curse of dimensionality*, statistically expressive $n$-gram models require enormous amounts of data. Indeed it can be expected that most training corpora will miss many plausible sequences, causing sparsity [Jurafsky and Martin, 2019, Ch. 3.4].

A substantial limitation is that the approach does not reuse information among similar words or word sequences. Hence, a unique representation is required per $n$-gram. The next section describes some approaches that use vectorized representations that exploit similarities between words to alleviate the dimensionality issue.

## 3.2 Word embeddings

The *distributional hypothesis* [e.g., Harris, 1954] states that words with similar meanings share similar contexts. So, if words tend to co-occur in the same sentences or documents, their meanings are assumed to be related. This assumption gives rise to *distributed representations*

of words, i.e. *word embeddings* [Jurafsky and Martin, 2019, Ch. 6]. These are fixed-length real-valued vectors. Distances between vectors correspond to similarities between words.

Creating rich representations of words and solving the language modeling task can be interpreted as two subproblems [Bengio et al., 2003]. Word2vec [Mikolov et al., 2013a] and its successor FastText [Bojanowski et al., 2017] are approaches to obtain word embeddings with neural networks efficiently.

### 3.2.1 Word2vec

Word2vec is a tool[1] for computing rich word embeddings. It implements the algorithms Continuous Bag-of-Words (CBOW) and Skip-gram [Mikolov et al., 2013a].



Figure 3.1: Schema of word2vec's CBOW and Skip-gram algorithms. Designed after Mikolov et al. [2013a].

**Continuous Bag-of-Words**

The CBOW model predicts a current word based on its context. It consists of an input, projection, and output layer (see Figure 3.1) [Mikolov et al., 2013a]. Several context words are given as input within a specified window (e.g., four future and four history words). The projection layer is shared for the input words (ignoring their order) and projects into a single position. The Bag-of-Words (BOW) model is similar in that it ignores the order of words. However, it produces sparse vectors from a term-document matrix. CBOW, as opposed to this, uses continuous distributed representations of the context.

**Skip-gram with negative sampling**

In contrast to CBOW, the Skip-gram model starts at a current word and tries to predict its context (see Figure 3.1). It does so by applying a log-linear classifier with a continuous projection

---

[1]https://code.google.com/archive/p/word2vec/

layer to each current word. Given a word sequence $w_1, w_2, ..., w_T$, the model is trained to maximize the following objective [Mikolov et al., 2013b]:

$$\frac{1}{T} \sum_{t=1}^{T} \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j}|w_t) \tag{3.4}$$

with $c$ denoting the size of the training context and $T$ the number of words. Mikolov et al. [2013b] approximate this using *negative sampling*. Its underlying idea is that a good model should be able to distinguish noise from data. Thus, the authors suggest to define $\log p(w_{t+j}|w_t)$ in Equation 3.4 as:

$$\log \sigma(v'_{w_O}{}^T v_{w_I}) + \sum_{i=1}^{k} \mathbb{E}_{w_i \sim P_n(w)}[\log \sigma(v'_{w_i}{}^T v_{w_I})] \tag{3.5}$$

where $\sigma(x) = 1/(1 + \exp(-x))$. $P_n(w)$ is a noise distribution from which $k$ negative samples are drawn. The objective is to differentiate those from the target word $w_O$ via logistic regression.

### 3.2.2 FastText

FastText [Bojanowski et al., 2017] is a continuation of word2vec [Mikolov et al., 2013a,1].

#### Incorporating subword information

German contains many different forms of one word, e.g., through verb inflections. Additionally, compound words, which bind together different nouns ("table tennis" becomes "Tischtennis", where "Tisch"="table" and "Tennis"="tennis", example from Bojanowski et al. [2017]) are commonly used. As compound words can be created this way freely, it is intractable to capture all of them in a model's vocabulary. FastText [Bojanowski et al., 2017] introduces subword information to Skip-gram with negative sampling [Mikolov et al., 2013b] (Section 3.2.1) to allow the sharing of representations across words. This adjustment helps with efficiency and with the handling of out-of-vocabulary (OOV) words. In Bojanowski et al. [2017] the approach outperformed word2vec on German text.

#### Modified Skip-gram

In Equation 3.2.1, a vector $v_{w_I}$ represents one specific input word. In FastText, a vector represents not a word but a *bag-of-character n-gram*. For a specified $n$, the word is split into all possible character $n$-grams (i.e., sequences of $n$ consecutive characters) and the complete word itself is also added as a special sequence to this set $G_w$. To distinguish between prefixes, suffixes, and other character sequences, boundary symbols $<$ and $>$ are added to the beginning

and end of words. The following is an example for the word *where* and $n = 3$ taken from Bojanowski et al. [2017][2]: *<wh, whe, her, ere, re>, <where>*

For each character sequence $g$ in $G_w$, a vector representation $z_g$ is computed and the sum of all $z_g$ forms the aggregated word representation. Applying these changes to Equation 3.5 yields the following new objective:

$$\log \sigma(\sum_{g \in G_w} {v'_{w_O}}^T z_g) + \sum_{i=1}^{k} \mathbb{E}_{w_i \sim P_n(w)}[\log \sigma(\sum_{g \in G_w} {v'_{w_i}}^T z_g)] \qquad (3.6)$$

OOV words can be represented by leaving out the special sequence (vector of the complete word) and by aggregating the matching $n$-gram representations [Bojanowski et al., 2017].[3]

## 3.3 Modeling sequences with neural networks

Sequential data is often modeled with types of recurrent neural networks (RNNs) [Rumelhart et al., 1986]. A more recent type of network that has gained widespread popularity is the Transformer network [Vaswani et al., 2017].

### 3.3.1 Recurrent neural networks

RNNs are a class of networks that reuse previous outputs as inputs and share weights across timesteps [Goodfellow et al., 2016, Ch. 10]. The following sections explain its basic form and two variants: Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber, 1997] and Gated Recurrent Units (GRU) networks.

**Vanilla recurrent neural networks**

The value of an RNN's hidden unit is called hidden state $h$. At timestep $t$ it is defined as:

$$h_t = g(h_{t-1}, x_t; \theta) \qquad (3.7)$$

where $g$ is an activation function, and $\theta$ model parameters that are shared over time [Goodfellow et al., 2016, Ch. 10.1]. The current input is denoted $x_t$.

The recurrent function $g$ can be unfolded over time. This property allows training via backpropagation equivalently to a feedforward structure. This is called *backpropagation through time* [Rumelhart et al., 1986].

Computation of the gradient of the first hidden state $h_1$ involves repeated multiplication with the derivative of the non-linear activation function $g$. This causes an exponential decrease

---

[2]The authors propose to extract all $n$-grams for $3 \leq n \leq 6$ [Bojanowski et al., 2017].

[3]Bojanowski et al. [2017] could show that using $n$-gram aggregates for OOV words led to improvements over the Skip-gram with negative sampling [Mikolov et al., 2013b] and CBOW [Mikolov et al., 2013a], as well as a version of FastText where null vectors represented OOV words.

or increase of weights, depending on the magnitude of the derivative. Vanishing or exploding gradients are a general problem of basic RNNs for long sequences [Hochreiter and Schmidhuber, 1997]. This phenomenon diminishes past inputs and complicates the learning of long-term dependencies of sequences longer than ten [Bengio et al., 1994].

**Long Short-Term Memory**

Hochreiter and Schmidhuber [1997] introduced the LSTM as a *gated* RNN version to deal with the problem of vanishing and exploding gradients. This model uses a linear cell state that facilitates gradient flow over longer sequences (*long-term memory*). Gates surrounding the cell control the flow of new and old information (*short-term memory*).

**The gates** The *forget gate* outputs a weight matrix $f$ with values between 0 and 1 to regulate how much of the previous cell state is discarded per entry, upon consideration of the past hidden state $h_{t-1}$ and current input $x_t$:[4]

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \tag{3.8}$$

where $W$ is the gate-specific weight matrix and $b$ the respective bias. The parameterization of the gates is learned [Hochreiter and Schmidhuber, 1997]. Similarly, the *input gate* controls which information is added (via weight matrix $i$) and the *output gate* controls which information to output (via $o$). Their definitions are equivalent to Equation 3.8.

**Controlling information flow** The equations below show how the gates are applied to determine the new cell and hidden states. The candidate values for the state update $z_t$ are scaled by $i_t$, and the past cell state $C_{t-1}$ is regulated by $f_t$ before entering the new cell state. Finally, the new hidden state $h_t$ is given by $C_t$ scaled by $o_t$.

$$z_t = tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \tag{3.9}$$

$$C_t = f_t \cdot C_{t-1} + i_t \cdot z_t \tag{3.10}$$

$$h_t = o_t \cdot tanh(C_t) \tag{3.11}$$

Information that is irrelevant (and thus discarded through the gates) does not get propagated back in time [Hochreiter and Schmidhuber, 1997]. Backpropagation happens only between cell states $C_t$ and $C_{t-1}$, without decay through non-linearity (since $C_{t-1}$ is only scaled through elementwise multiplication with $f_t$). Signals within the cell state "can flow back indefinitely without ever being scaled" [Hochreiter and Schmidhuber, 1997, p. 1747]. This alleviates the problem of vanishing or exploding gradients.

---

[4]Notations were taken from `https://colah.github.io/posts/2015-08-Understanding-LSTMs/`.

**Gated Recurrent Units**

GRUs [Cho et al., 2014a] are a modified type of LSTM [Hochreiter and Schmidhuber, 1997] units, which do not distinguish between a cell memory and the hidden state. Furthermore, the gates are merged and reduced to a *reset* and an *update gate* to scale down the number of parameters. The update gate combines the functions of the LSTM's forget and input gates. Its weights determine the extent of the unit's activation update. The reset gate controls how much of previous hidden states is forgotten. So, just like LSTMs, GRUs allow maintaining previous information while adding new information, such that essential features are not fully overwritten [Chung et al., 2014]. The additive nature of the model creates shortcuts for the gradient flow and prevents vanishing gradients. The simplifications do not per se cause performance loss. Depending on the task, GRUs were shown even to outperform LSTMs [Chung et al., 2014].

**Learning the probability distribution over a sequence**

Training an RNN-based language model can be achieved by letting it learn to predict the next word in a sequence. At each timestep, the most probable word (represented as a vector) can be computed via a softmax over the vocabulary [Cho et al., 2014b]. The product of its conditional probabilities gives the probability of a sequence (see Equation 3.2). Sampling a new sequence with the trained model is done by iteratively sampling a word per timestep.

### 3.3.2 RNN encoder-decoder

In encoder-decoder architectures, the encoder function maps the input to a latent representation from which the decoder projects to the target. This type of architecture can be used to solve *sequence-to-sequence* tasks, like neural machine translation. Respective approaches introduced by Cho et al. [2014b] and Sutskever et al. [2014] use two RNN-type networks, where one serves as the encoder and the other as the decoder. The variable-length input sequence is encoded to a fixed-length vector representation and then decoded to, again, to a variable-length output.

The networks are jointly trained to model the conditional distribution of a sequence on another sequence, with lengths $T$ and $T'$ (these can differ):

$$p(y_1, ..., y_T | x_1, ..., y_{T'}) \tag{3.12}$$

[Cho et al., 2014b]. The encoder RNN processes the entire input sequence sequentially and creates a summary in the form of a fixed-size *context vector* $c = q_1(h_1, ..., h_T)$ ($q$ is an activation function) [notation from Bahdanau et al., 2014]. The decoder RNN then generates the output by predicting the next word based on the past hidden state $h_{t-1}$, its own previous generation $y_{t-1}$, as well as $c$. The consumption of previously generated output for generation of the next output is termed *auto-regression* [Vaswani et al., 2017].

The current decoder hidden state is given by Equation 3.13 [Cho et al., 2014b] and the next symbol's conditional probability is defined in Equation 3.14. Both $f$ and $g$ are activation functions, while $g$ returns probabilities, for example, via softmax [Cho et al., 2014b].

$$h_t = f(h_{t-1}, y_{t-1}, c) \tag{3.13}$$
$$P(y_t | y_{t-1}, y_{t-2}, ..., y_1, c) = g(h_t, y_{t-1}, c) \tag{3.14}$$

**Introducing attention**

Squashing long sequences to fixed-sized representations can cause information loss and again complicates the modeling of long-range dependencies. Bahdanau et al. [2014] first used *attention* mechanisms in an RNN-based encoder-decoder to emphasize the most relevant information of an input sequence when generating the output sequence. Instead of a single summarizing vector, their approach yields a context vector for each target word to alleviate the loss of valuable information.

A context vector summarizes information of the full input sequence. Through a weighting mechanism, context vector $c_t$ (at timestep $t$) puts emphasis on input words that align with the output position $t$. The alignment is determined by a trained *alignment model* [Bahdanau et al., 2014]. This allows the decoder to pay attention to relevant information. Furthermore, it loosens up the strict sequentiality in that information is shared throughout the sequence [Bahdanau et al., 2014].

### 3.3.3 Transformer network

In general, RNNs are sequential and, thus, do not allow for parallelization. This poses limitations on the modeling of long sequences due to memory constraints [Vaswani et al., 2017]. With the Transformer model, Vaswani et al. [2017] presented an approach that does not rely on recurrence but on attention only. With this, it fully exploits the effect of information flow throughout a sequence (Section 3.3.2).

Word2vec yields one vector representation for one word, independent of its current context or position. Transformer-based models, in contrast, are sensitive to this information due to the emphasis on attention, i.e., they are *contextualized*.

**Architecture**

The Transformer is a neural network that also follows an *encoder-decoder* structure. The encoder is a stack of six encoder layers, and the decoder a stack of six decoder layers (see Figure 3.2). However, information is not processed in a strictly sequential manner. Instead, an entire sequence (within a certain window size) is encoded and fed into the decoder. The decoder additionally receives its previous outputs in an auto-regressive fashion [Vaswani et al., 2017].

Figure 3.2: Schematic overview of the Transformer encoder-decoder architecture. Designed after Vaswani et al. [2017].

## Adding positional information

The network input is a sequence of words embedded into a vector of a fixed size $d_{model} = 512$. Vaswani et al. [2017] use positional encodings to infuse information about the sequence order. The encodings are vectors of dimensionality $d_{model}$ that are added to the embeddings to adjust the distances between the word representations.

## Self-attention

A defining characteristic of the Transformer is its use of *self-attention*. It allows representing a sequence via the relationships of words within this sequence. This way, a current word is processed concerning its surrounding words, i.e., contextualized.

A self-attention function gets as input *queries* and sets of *key-value pairs*. The values carry the content information of the input, and the keys serve as an indexing mechanism. The queries determine which information is accessed. This is done via a compatibility metric (e.g., the dot-product) between key and query.

For parallelization, Vaswani et al. [2017] combine multiple queries, keys, and values in matrices $Q$, $K$, and $V$. These are derived from the word embeddings through learned projection matrices. Finally, the self-attention function is a *scaled dot-product* formalized as:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V \tag{3.15}$$

where $d_k$ is the dimensionality of query and key vectors, and the results are scaled by $\sqrt{d_k}$. A softmax function is applied to normalize the resulting scores, which then serve as weights.

**Multiple attention heads**  Vaswani et al. [2017] use eight different, randomly initialized projection matrices to allow attending to information from different representation subspaces at various positions. These different attention layers are called *heads*. The approach is, hence, called *multi-head attention*. The decoder's self-attention sub-layer uses masking to ensure that the predictions for the current position only depend on previous predictions (see Figure 3.2, *masked* multi-head attention, marked **b.**).

As can be seen in Figure 3.2, the encoder and decoder communicate via a multi-head attention function (marked **c.**). This operation considers the keys and values from the encoder stack output. This interface allows the decoder to attend to the currently relevant input.

## Network output

A linear layer follows the decoder to produce a logits vector from the decoder stack output. The size of the logits vector equals the model's vocabulary size and provides a score for each word in the vocabulary. Finally, softmax is applied to interpret the scores as relative probabilities.

## 3.4 Language models with contextualized embeddings

The following sections present two types of Transformer architectures: BERT [Devlin et al., 2019] is an encoder model that yields contextualized language representations. GPT-2 [Radford et al., 2019] and GPT-3 [Brown et al., 2020] are decoder-based and generative models.

### 3.4.1 BERT

After pre-training on large corpora of text data, BERT (Pre-training of Deep Bidirectional Transformers for Language Understanding) [Devlin et al., 2019] yields language representations that can be used for different downstream tasks via transfer learning.

**Architecture**

BERT is a Transformer network with an architecture similar to the Transformer **encoder** stack introduced in Section 3.3.3 but more highly parameterized at different points of the architecture.[5]

**Pre-training**

Devlin et al. [2019] used two unsupervised training objectives: *masked language modeling* (MLM) and *next sentence prediction* (NSP). For MLM, one of the input tokens is randomly masked with a [MASK] token, and the original vocabulary ID is predicted based on the remaining tokens. The objective is a cross-entropy loss on the prediction. These tokens represent the left and right context of the masked token, yielding *bi-directionality*.

For the NSP task, two token segments are fed into the model and separated and followed by a [SEP] token. The second segment corresponds to a randomly chosen sentence, e.g., from another document or the actual successive sentence. The objective is a binary classification loss for predicting if the second segment is the true successor.

**Classification**

A special [CLS] token (*CLS* for classification) is prepended to each input sequence. After processing the sequence, the output vector corresponding to the [CLS] token carries information about the complete sequence. Feeding this vector to a feedforward network provides a classifier.

However, other ways were introduced to improve classification performance further. One alternative to get a fixed-size sentence representation is to average all word embeddings from the output layer [Reimers and Gurevych, 2019]. Reimers and Gurevych [2019] further improved the semantic meaningfulness of sentence embeddings through the finetuning of assembled BERT networks with shared weights. The approach is described in more detail in Section 3.4.1.

---

[5]BERTBase comprises 12 encoder layers, has 768 hidden units per feedforward network, and 12 attention heads for the multi-head attention layers. BERTLarge has 24 encoder layers, 1024 hidden units per feedforward network, and 16 attention heads.

### RoBERTa

RoBERTa is an improved version of BERT that uses the same architecture but different objectives. Liu et al. [2019] changed the MLM objective to use a dynamic masking pattern. They further removed the NSP objective and instead trained the model on longer input sequences. The authors also trained RoBERTa on a significantly more extensive data set and for more training passes.[6]

For BERT, Devlin et al. [2019] used WordPiece tokenization to embed words based on a vocabulary of size 30K. WordPiece is a subword-level encoding strategy that uses characters as a base vocabulary and merges them by trained rules. Vaswani et al. [2017] call this *character-level* byte-pair encoding (BPE). RoBERTa, in comparison, utilizes *byte-level* BPE, where bytes serve as the base vocabulary. This facilitates better handling of vast data sets and allows for a larger vocabulary (50K entries). RoBERTa outperforms BERT on several NLP benchmarks mainly due to increased training data and duration [Liu et al., 2019].

### SentenceBERT

As explained in Section 3.4.1, a common strategy to retrieve sentence representations from BERT and RoBERTa is to either average the output layer (MEAN strategy) or to take the [CLS] token embedding (CLS strategy). Reimers and Gurevych [2019] claim that these methods do not yield satisfying semantics on a sentence level. Their contribution, SentenceBERT, performed better on Semantic Textual Similarity (STS) benchmarks [e.g., Agirre et al., 2016; Cer et al., 2017], which indicates improved semantic meaningfulness.

**Pairing networks**　SentenceBERT [Reimers and Gurevych, 2019] is based on pre-trained BERT or RoBERTa embeddings. Multiple networks (two or three, depending on the task) are combined via shared weights (see Figures 3.3). The resulting *siamese* or *triplet network* is then finetuned towards a sentence-based objective. Each of the original networks receives one sentence as input. Their outputs are pooled with the MEAN strategy mentioned above.[7]

**Classification task**　For a classification task, the pooled sentence embeddings, $u$ and $v$, are concatenated with each other and their element-wise difference $|u - v|$ (see Figure 3.3a for a schematic overview). This difference term is a distance measure and ensures similar sentence pairs are represented closer than dissimilar ones. The concatenated embeddings are then multiplied with a trainable weight matrix and followed by a softmax classifier.

---

[6]RoBERTa was trained for 500K steps and a batch size of 8K sequences on 160 GB of data [Liu et al., 2019].

[7]Other strategies, like the CLS strategy or a max-over-time strategy (to capture the essential features in a sentence), exist but were shown to work less well [Reimers and Gurevych, 2019].

(a) Siamese network architecture for a classification task.

(b) Siamese network architecture for a regression task.

Figure 3.3: Siamese model architectures. Schema designed closely after Reimers and Gurevych [2019]. Each of these architectures uses a pair of either BERT or RoBERTa embeddings. Which of the objectives is used depends on the task at hand.

**Regression task**   For the regressor, the cosine-similarity between the embeddings $u$ and $v$ is computed, and the model is trained via mean-squared error (see Figure 3.3b). Because distances between sentence pairs are at the focus of the SentenceBERT tuning procedure, the resulting embeddings permit the computation of meaningful cosine-similarities.

### Multilingual embeddings

The original version of BERT and RoBERTa were trained in English. Obtaining embeddings in a different language from pre-trained versions requires extra steps. Multilingual language models own a shared embedding space for distinct languages [Lample and Conneau, 2019]. Reimers and Gurevych [2020] introduced an approach that takes a monolingual model, like SentenceBERT, as a basis for the training of a multilingual version.

**Teacher and student approach**   The approach uses a pre-trained teacher model $M$ for language $s$. It further requires a parallel corpus of sentences $s_i$ and their translations $t_i$. The translations can belong to different languages. Via mean-squared loss, a student model $\hat{M}$ is trained such that $\hat{M}(s_i) \approx M(s_i)$ and $\hat{M}(t_i) \approx M(s_i)$. So, the objective with respect to the student model is to maintain the original representations and to map the translations to the same location in the embedding space as their original counterpart.

### 3.4.2 GPT-2 and GPT-3

GPT-2 [Radford et al., 2019] and GPT-3 [Brown et al., 2020] are the successors of GPT (Generative Pretrained Transformer) [Radford et al., 2018]. These language models are also pre-trained on vast text corpora in an unsupervised fashion. They are trained on the standard language modeling task (Equation 3.1). The GPT-based models have similar architecture and differ mainly in parameterization.

#### Architecture

Whereas BERT's architecture is based on an encoder stack, the versions of GPT consist of Transformer **decoder** blocks (also introduced in Section 3.3.3, Vaswani et al. [2017]). The input to the network are *context* tokens (the start token <|endoftext|>, a generation prompt, or previously generated output).

**Modifications of the original decoder** GPT-2 introduced some modifications to the original decoder block. For example, layer normalization was moved to the input of each sub-block, and another layer normalization was added after the last self-attention block [Radford et al., 2019]. The context window was extended to 1024 tokens (from 512). For GPT-3, even larger context windows of 2048 were used. The models use byte-level BPE with an adjustment to prevent merging across character categories (e.g., combining letters and punctuation).

**High parameterization** Both GPT-2 and GPT-3 were pre-trained for different model sizes. The GPT-2 variants range from a 117M parameter model with 12 layers and $d_{model} = 768$ to a 1542M parameter version with 48 layers and $d_{model} = 1600$ [Radford et al., 2019]. For GPT-3, the smallest version has 125M parameters, 12 layers, $d_{model} = 768$ and the largest has 175B parameters, 96 layers, and $d_{model} = 12288$ [Brown et al., 2020].

#### Competent generalists

By providing the models with extensive amounts of data on various domains, the authors hoped to move past smaller, task-specific networks and to create "competent generalists" [Radford et al., 2019, p. 1] instead. The two networks are trained to qualify for *zero-shot* application to downstream tasks. This means that no architectural modification or training is needed before application.

# Chapter 4

# Data Collection for Regard Classification

The concept of *regard* was introduced by Sheng et al. [2019] for the specific purpose of measuring social bias in generative models. According to the authors' definition, an utterance of positive *regard* causes most people to think highly of a described person, and a description conveying negative *regard* causes others to think lowly of a person. If a language model's generated sentences systematically convey better *regard* towards one over another demographic, then the model perpetuates bias.

In order to make the concept quantifiable in the German context, a dataset containing German personal descriptions was collected and annotated. This dataset was later used for the training of a *regard* classifier (see Chapter 5), which can serve as a global bias measure [Liang et al., 2021]. Each description is a single sentence. These were crowd-sourced on an online platform, with participants recruited via social media. Section 4.1 illustrates the design of the questionnaire. Essential insights on the collected data are described in Section 4.2.

Section 4.3 describes how independent annotator ratings were collected to improve the data quality [Bernstein et al., 2010]. The collected annotations were evaluated through an inter-rater agreement metric. Two different ways to aggregate the multiple labels were explored: through a majority vote (*mode*) versus through a unanimous vote (*consensus*). Fortunately, the amount and nature of the data enabled counterfactual balancing to prevent gender bias [Lu et al., 2020; Maudslay et al., 2019], which is described in Section 4.4. This chapter is concluded by summarizing remarks on the data preparation steps.

## 4.1 Design of the crowd-sourced data collection

The goal of the data collection was to gather a diverse set of sentences, in which a person is described to be *regarded* as negative or positive, or in a neutral way. These valences form the three target classes. Each participant was assigned to one of the conditions so that they had to invent descriptions for one of the target categories exclusively. This was to maintain control over the number of participants per condition.

The survey was created with the software SoSci Survey [Leiner, 2021] and made available to users via the online platform www.soscisurvey.de. It was conducted in German and designed to collect up to 20 sentences per person. Participants entered each sentence into a text field preceded by a predefined prefix (later also referred to as *prompt*), as can be seen in Figure 4.1. The prefix consisted of at least "Die Person"/ "The person", and for the first ten items, it contained an additional context term (described in Section 4.1.1). With each prefix, the survey displayed a unique portrait of a person to inspire diverse descriptions. See Section 4.1.2 for a more detailed explanation.

### 4.1.1 Textual prompts

For the first ten items, the prompt-input pattern was (prompt in bold):
**["Die Person"]** + **[context term]** + [generated *regard* term] and for the final ten items:
**["Die Person"]** + [generated *regard* term]. The context terms were taken and translated from Sheng et al. [2019], as compared in Table 4.1. The authors distinguished between *respect* and *occupation* contexts. *Occupation* is a well-studied bias context in existing research [e.g., Bolukbasi et al., 2016; Rudinger et al., 2018]. The *respect* contexts complement these by adding a more descriptive expression of social bias [Sheng et al., 2019]. In this thesis, the ten no-context prompts were added to also encourage sentences beyond these contexts.

| This work | Sheng et al. [2019] | Bias context |
|---|---|---|
| ... arbeitete als | ... worked as | |
| ... hatte eine Stelle angefangen als | ... started working as | |
| ... verdiente Geld, indem | ... earned money by | Occupation |
| ... hatte eine Teilzeitstelle als | ... had a part-time job as | |
| ... hatte einen Job als | ... had a job as | |
| ... wurde wahrgenommen als | ... was thought of as | |
| ... galt als | ... was known for | |
| ... wurde betrachtet als | ... was regarded as | Respect |
| ... wurde beschrieben als | ... was described as | |
| ... war sehr bekannt dafür | ... was well-known for | |

Table 4.1: Bias contexts taken and translated into German from Sheng et al. [2019].

### 4.1.2 Visual cues

Each prompt was accompanied by the photo of an actual person, as shown in Figure 4.1. These were sampled from the FairFace database [Kärkkäinen and Joo, 2021].[12] The survey randomly sourced its images from pre-sampled subset of the corpus. This subset contained 168 images of people of color and 167 images of white people, evenly distributed across age ($> 9$) and gender.

### Insights from the piloting phase

In a dedicated piloting phase, feedback on preliminary versions of the survey was provided by ten pilot testers. This helped to avoid technical mistakes and test the usability of the programmed survey. Initially, the portrait photos were sampled evenly across ethnicity, yielding a set of majority ($\approx 86\%$) non-white people. Four out of the ten pilot testers reported feeling uncomfortable filling out the text fields. The presence of mainly non-white people in the images created the impression that the study was supposed to "trick" them or figure out their unconscious biases. One participant said that they put in extra effort to say nice things (positive condition). Another one said they felt awful about falling back to stereotypical and racist utterances (negative condition). These reports indicated the presence of a *social desirability* effect. This psychological construct describes the tendency to respond to questionnaires in ways that comply with the interviewer's expectations or social norms in general [Bogner and Landrock, 2016]. Although there were no correct or incorrect answers to the items, the perceived discomfort of the participants was reason enough to reconsider the use of the visual cues.

### Dealing with the social desirability issue

In response to the above issues, 50% of the photos were sampled to show people of color and 50% to show white people. Additionally, the instruction texts were adjusted to emphasize that personal opinions were not of interest. Participants were encouraged to invent content that transports the assigned type of *regard* properly, irrespective of political correctness. The role of the presented photos was further clarified as serving only as an inspiration. In the main collection phase, no more comments regarding this issue were received.

---

[1] https://github.com/joojs/fairface
[2] The FairFace dataset contains photos equally distributed for the following ethnicities: Middle Eastern, Indian, South-East Asian, East Asian, Latinx, Black, White [Kärkkäinen and Joo, 2021].

Vervollständigen Sie den Satz mit einer POSITIVEN Beschreibung. Die Länge der Sätze darf variieren. Sie können kurze, aber auch sehr lange Sätze schreiben.

Sollten Sie das Geschlecht der Person nicht genau ausmachen können, wählen Sie frei aus. Denken Sie daran, das Bild soll nur Ihre Fantasie anregen. Ihre Beschreibung muss nicht zwangsläufig auf die gezeigte Person passen.

| Die Person wurde wahrgenommen als | |
|---|---|

(a) Prefix: *"The person was thought of as"*



Vervollständigen Sie den Satz mit einer NEUTRALEN Beschreibung. Die Länge der Sätze darf variieren. Sie können kurze, aber auch sehr lange Sätze schreiben.

Sollten Sie das Geschlecht der Person nicht genau ausmachen können, wählen Sie frei aus. Denken Sie daran, das Bild soll nur Ihre Fantasie anregen. Ihre Beschreibung muss nicht zwangsläufig auf die gezeigte Person passen.

Denken Sie sich einen eigenen Satz aus – dieses Mal ohne vorgegebenen Satzanfang.

| Die Person | |
|---|---|

(b) Prefix: *"The person"*

Figure 4.1: Example questionnaire pages from the data collection. The presented photos varied and were presented to inspire diverse personal descriptions. With reference to Section 4.1.1: (a) Uses a contextualized prompt. (b) Uses an open prompt.
Translated instruction for the screenshot (a): *"Complete the sentence with a POSITIVE description. The length of the sentences may vary. You can write short, but also very long sentences. If you cannot exactly make out the gender of the person, choose freely. Remember, the picture is only to stimulate your imagination. Your description does not necessarily have to fit the person shown."*

45

## 4.2 Collected data

This section summarizes some insights on the final participants of the study and data quality.

### 4.2.1 Participants

Participants were recruited within the company-internal messenger, through the personal academic network, and via the professional networking platform LinkedIn.[3] Participants had to explicitly give consent, be above 18, and have native-level German skills.[4]

**Demographics**

Unfortunately, detailed demographic information was not collected from the very beginning. However, this type of information can be crucial to understand how the content of a dataset might be biased [Bender and Friedman, 2018] - especially for this very subjective task of describing people in evaluative ways.

Based on the incomplete data collected via the survey and additional information drawn from the social media accounts of people that publicly indicated participation, the estimated ratio of male participants is $2/3$. Most participants were recruited within a network of mainly white, well-educated persons. The majority were actively working or studying at the university level, and many were associated with a technological field, marketing/sales, or design. This estimation should raise the awareness that the data was sourced from a crowd with limited diversity. This can cause bias and decrease generalizability [Bender and Friedman, 2018].

**Response rates**

Out of 116 people who launched the survey, 75 provided at least one sentence, and 50 filled out the entire questionnaire. Finishing the complete questionnaire took between 5 and 22 minutes ($\approx 12$ minutes on average). Participants skipped $31\%$ of all items in the negative condition and $36\%$ of items in the positive condition. In the neutral condition, only $17\%$ were left unresponded. Consequently, more neutral sentences were generated than negative or positive ones (Figure 4.2).

Since the data was collected anonymously, it was impossible to ask specific participants why they dropped out or skipped on questions. However, some of the general feedback indicated that it was challenging to disparage or only praise an invented person. The set of neutral responses, on the other hand, contains many descriptions of what a person does or wears, which may be easier to imagine.

---

[3] https://www.linkedin.com/

[4] If a participant answered one of the corresponding queries with "No", they were automatically redirected to the last page of the online survey (8 persons in total).

Figure 4.2: Pre-annotation label distribution: Number of sentences by survey condition. A total of 1,157 sentences was collected. The left plot includes prefixes with and without bias context.

### 4.2.2 Descriptive statistics and content exploration

A total number of $1,157$ sentences was collected. Figure 4.2 shows the distribution across conditions. The dataset contains $2,599$ unique words ($2,457$ without stopwords[5] and $4,447$ unique bigrams. The average amount of characters in a sentence is $38.1$ ($SD = 26.2$, $range = [2, 203]$) and the average count of words is $5.5$ ($SD = 4.1$, $range = [1, 34]$).

The maximum frequency for the non-stopwords is $35$. Only the following words appeared more than $20$ times: "immer"/ "always", "gerne"/ "with pleasure", "Menschen"/ "humans", and "wurde"/ "was" or "became".

#### Qualitative description

The word frequencies shown in Figure 4.3 intend to give an intuition of the sentence contents. The plots show the most frequent lemmatized nouns and adjectives.[6] The negative condition contains many expected words, like "schlecht"/ "bad", "Droge"/ "drug", "Prostituierte"/"prostitute". Note that "Schlange"/"snake" was used to characterize someone as deceitful.

For the positive condition, terms like "gut"/"good", "Erfolg"/"success", and "Lachen"/ "laughter" are expected, as well. Finally, the neutral condition contains many profession-related mentions, like "Lehrerin"/ "teacher", "Schule"/ "school", "Arbeit"/"work", and "Unternehmen"/ "business". Though, beyond that, the neutral and positive conditions show much overlap. For both, "gut"/ "good" and "Leben"/"life" appears amongst the highest ranks. This observation motivated an independent annotation step described in the following Section 4.3.

---

[5] The German stopword list of the Natural Language Toolkit (NLTK) (https://www.nltk.org/) containing 232 words was used.

[6] The following frequent words were removed from these lists due to their overlap across conditions and little informativeness: "Mensch"/"human", "Freund"/"friend", "Frau"/"woman", "Mann"/"man", "Kind"/"child", "Mitmensch"/"fellow human being", "Leute"/"people", "Person"/"person".

Figure 4.3: Lemmatized top-15 most frequent nouns and adjectives per questionnaire condition.

## 4.3 Crowd-sourced annotation

Crowd-sourcing data through an anonymized procedure with little interviewer-interviewee interaction comes at the cost of "uncertain worker quality" [Bernstein et al., 2010, p. 316]. Following up the crowd-sourced data collection with an independently crowd-sourced verification step can reduce noise and improve the output quality [Bernstein et al., 2010]. Hence, to ensure that the collected sentences are meaningful, understandable, and consistent with the intended *regard* concept, an independent verification and labeling step was performed.

Please note, the annotation procedure presented below differs slightly from the original by Sheng et al. [2020]. The two procedures were compared in an additional experiment, which is reported in Appendix A.2. The results indicate that the variant used here is an improvement on the original.

### 4.3.1 Independent annotator ratings

One female and four male students (fellow students in the Master's course) helped out as annotators. Each received a written instruction on the concept of *regard* and labeling rules (see Appendix A.1). The annotators had to label each sentence of the previously crowd-sourced dataset as $-1$ for negative, $0$ for neutral, and $1$ for positive. An additional labeling option named *not sure* allowed the marking of sentences that were hard to understand or did not appear to fit into any of the categories. This fallback option introduced a verification aspect to the procedure. Later on, sentences with at least one *not sure* marking were dropped (23 cases). Only three sentences were marked as such by more than one person.

In 65% of the cases (752 of the 1,157 sentences), the annotations were unanimous, indicating a clear *regard* signal in most of the data set. For only 21 of the remaining 405 sentences, the annotators used the full range of possible labels, hinting at high ambiguity. Some of these ambiguous sentences are listed in Table A.3, in Appendix B.3.1.

### 4.3.2 Inter-rater reliability

Measures of *inter-rater agreement* provide a criterion for determining consistency among multiple raters [Gwet, 2008]. *Regard* ratings are inherently subjective. An example that illustrates the subjective nature of the labeling task is the following statement: "The person was a follower of Greta Thunberg". Whether this statement is considered praise or insult depends on one's political viewpoint. Hence, a perfect agreement between annotators is not expected. Still, a *good* level of the agreement should nevertheless be present to ensure that the conceptual understanding is coherent. A sentence like "The person was described as a triple murderer." should produce the negative *regard* label consistently across annotators.

## Cohen's kappa

Table 4.3 shows *Cohen's kappa* [Cohen, 1960] inter-rater agreement scores across each annotator. Cohen's kappa is a two-rater metric defined as

$$\kappa_{Cohen} = \frac{p_a - p_e}{1 - p_e} \tag{4.1}$$

where $p_a$ is the overall agreement probability and $p_e$ the probability of agreement by chance [Gwet, 2008]. The idea is that of *"contrasting observed with expected agreement"* [Fleiss et al., 2003, p. 602]. Given two annotators $A$ and $B$, the probability $p_a$ is given by:

$$p_a = \sum_{k=1}^{q} p_{kk} \tag{4.2}$$

where $p_{kk} = n_{kk}/n$, with $n$ denoting the amount of data points to be labeled, and $n_{kl}$ the number of data points that $A$ labeled as $k$ and the $B$ labeled as $l$ [Gwet, 2008]. Further, $q$ represents the number of target classes. The chance-agreement $p_e$ is defined as:

$$p_e = \sum_{k=1}^{q} p_{Ak} p_{Bk} \tag{4.3}$$

where $p_{Ak} = n_{Ak}/n$, the proportion of sentences that $A$ labeled as $k$ (equivalently so for $p_{Bk}$) [Gwet, 2008].

## Levels of agreement between annotators

The Cohen's kappa values among Annotator 0 to Annotator 4 in Table 4.3 range from *moderate* to *strong*, with reference to the norms introduced by McHugh [2012] (see Table 4.2).[7] With an average of .80, sufficient data reliability is indicated.

McHugh [2012] emphasizes that agreement below 1.0 is also a measure of *disagreement* among raters and an indicator of data (un-)reliability. In the clinical context, rater disagreement indicates that the data misrepresents the research object since there is only one truth [McHugh, 2012]. In the context of subjective concepts like *regard*, however, this can behave differently. Whenever annotators disagree, there may be more than one correct label. This is, for example, the case with political or religious statements. In these cases, we cannot hope for a clear *regard* signal, and it might, moreover, not make much sense to try to define a norm. Nevertheless, it appears helpful to keep in mind that low levels of agreement can indicate that the data diverges from the construct that is to be measured.

---

[7]Values below 0 are possible and should be interpreted as *no agreement*. However, the authors argue that they are not meaningful and likely caused by mistakes in the data collection.

| $\kappa$ | Level of agreement |
|---|---|
| 0 - .20 | None |
| .21 - .39 | Minimal |
| .40 - .59 | Weak |
| .60 - .79 | Moderate |
| .80 - .90 | Strong |
| > .90 | Almost perfect |

Table 4.2: Interpretation guide for Cohen's kappa adapted from McHugh [2012].

| | Ann. 0 | Ann. 1 | Ann. 2 | Ann. 3 | Ann. 4 | Original |
|---|---|---|---|---|---|---|
| Study conductor | .77 | .69 | .74 | .83 | .83 | .60 |
| Annotator 0 | | .75 | .73 | .80 | .80 | .59 |
| Annotator 1 | | | .74 | .71 | .71 | .56 |
| Annotator 2 | | | | .73 | .74 | .58 |
| Annotator 3 | | | | | .83 | .57 |
| Annotator 4 | | | | | | .60 |

Table 4.3: Cohen's kappa scores expressing inter-rater reliability between all annotators, the study conductor, and the original labels derived from the survey conditions. Values between .61 and .80 indicate *good agreement* [McHugh, 2012].

**Agreement between annotators and study conductor**

Table 4.3 also shows how well the annotators' understanding of *regard* was in line with the study conductor. With an average Cohen's kappa of .77, the independent annotators and interviewer ratings agree with each other *moderately*, and the constructs appear to match within expected levels. Substantial divergence would indicate that the instruction materials convey a different idea from the one intended. For the final dataset, the interviewer annotations were ignored to avoid bias.

**Agreement between annotators and study condition**

The qualitative data exploration in Section 4.2.2 raised the question of how well the collected sentences matched their study condition. As Table 4.3 illustrates, the average Cohen's kappa between the independent annotators and the original condition was .58. Since the annotators agree strongly with each other but weakly with the original conditions, it may be inferred that the original labels are less reliable. This supports the assumption that the separation of data creation and verification can help to denoise [Bernstein et al., 2010]. The original labels are henceforth ignored.

### 4.3.3 Aggregation of individual annotations

It appeared attractive to use disagreement as a criterion to identify items that do not or not only fit into one of the categories negative, neutral, or positive *regard*. Reducing the dataset to only items with full consensus was used to further denoise the data and distill it into sentences with a clearer signal towards the target classes.

Additionally, majority-based labeling was used as well (as was done in Sheng et al. [2020]). The majority opinion on the *regard* conveyed in a sentence could be interpreted as the quantification of a social norm. In summary, two different aggregation methods were experimented with to find which yields better performance of the *regard* classifier:

- **Mode annotation**: In taking the mode of all annotator ratings, majority labels for all of the sentences were computed (excluding the 23 marked *not sure* that were dropped earlier).

- **Consensus annotation**: As mentioned in Section 4.3.1, all annotators fully agreed on the valence of 65% of the sentences. So, as a second aggregate labeling method, these sentences were extracted from the corpus and assigned the corresponding consensus label.

Figures 4.4 and 4.5 show the distributions for both annotation types. For prompts with bias contexts, the distributions differ. Occupation-related sentences are mostly rated as neutral, respect-related sentences are more polar.



Figure 4.4: Frequencies of *regard* labels computed via the mode of all annotator ratings. The distributions for occupation and respect contexts differ. Note that the leftmost plot includes prefixes with and without bias context.

Figure 4.5: Frequencies of *regard* labels with full annotator consensus. The distributions for occupation and respect contexts are similar to Figure 4.4.

## 4.4 Gender-balancing the dataset

Counterfactual balancing of the training data can prevent the trained model from capturing this dimension without affecting the task performance [Barikeri et al., 2021; Lu et al., 2020]. Hence, a CDS-inspired [Maudslay et al., 2019] (explained in Section 2.2.1) procedure was applied to introduce gendered prefixes in a balanced way.

### 4.4.1 Distribution of gender indications

Explorative analysis of the data revealed that the survey's differently gendered visual cues (Section 4.1.2) yielded diverse use of gendered pronouns. Due to the relatively small sample size, all sentences could be manually labeled as male, female, or *none* if there was no gender indication. Figure 4.6 shows that male and female mentions were about evenly distributed. Figure 4.7 illustrates gender distributions across *regard* labels.



Figure 4.6: Distribution of subject genders as naturally occurring in the collected sentences. Most sentences did not indicate a gender (*none*).

### 4.4.2 Introducing gender-balanced prefixes

The consistent sentence structure (["Die Person"] + [context term] + [generated *regard* term]) made it easy to replace the subject with counterfactual demographic mentions ("Der Mann"/"The man" and "Die Frau"/"The woman"). Knowing which sentences contain gender markings allowed to maintain grammatical soundness. For the remaining gender-neutral sentences, the demographic mentions were assigned randomly with a 50% chance.



Figure 4.7: Distribution of labels by occurrence of subject gender in the sentence.

## 4.5 Concluding remarks on the dataset creation

A *regard* dataset with 1,157 sentences was collected. In a separate procedure, the ratings of independent annotators were used to denoise the data and create coherent labels for the development of the *regard* classifier. The agreement between annotators was high, indicating a reliable *regard* signal.

Both, the crowd-sourced *mode* and *consensus* labels were later used as development and testing targets during experimentation. As the results in Section 5.5 indicate, the added variability through the *mode* targets were beneficial to the training.

Gendered subjects were introduced to the dataset through the balanced insertion of gender-specific prefixes, following the idea of the CDS method for bias mitigation [Maudslay et al., 2019]. The evaluation in Section 5.6.3 confirms that the resulting *regard* classifier does not show signs of a gender bias.

# Chapter 5

# Development of a Regard Classifier

One way to evaluate if a language model is socially biased is to analyze its generated samples. We can infer a social bias if these outputs transport a skewed representation of certain demographic groups [Crawford, 2017].

A dedicated *regard* classifier can identify the global bias expressed in open-ended text [Sheng et al., 2019]. It can evaluate large amounts of samples with little effort to obtain comparable score ratios for different demographics. For this reason, the dataset described in Chapter 4 was used to train a German *regard* classifier.

As a baseline model, a Gradient Boosted Trees (GBT) [Friedman, 2001] model was used (Sections 5.1). TF-IDF-weighted FastText [Bojanowski et al., 2017] embeddings averaged for the words in a sentence were provided as input. Next, it was tested how sequentiality introduced by Gated Recurrent Units (GRU) [Cho et al., 2014a] would improve the performance (Section 5.2). For this, FastText embedded word sequences were fed to a GRU model. The final approach, a SentenceBERT-based [Reimers and Gurevych, 2019] classifier (Section 5.3), outperformed the previous two (Section 5.5). Given the short, single sentence descriptions, the information added through the contextualized SentenceBERT embeddings appeared especially valuable. The final model was thoroughly evaluated and tested for potential biases (Section 5.6).

## 5.1 Baseline model: FastText and Gradient Boosted Trees

The different approaches in this thesis utilized two different types of pre-trained word or sentence embeddings to vectorize the input text. The first two approaches were based on FastText [Bojanowski et al., 2017], which is explained in detail in Section 3.2.2.

### 5.1.1 Input preparation

For the baseline model, FastText embeddings per word in a sentence were weighted with their *inverse document frequencies* (explained below) to contrast descriptive from non-descriptive words. The weighted word vectors were then averaged and fed into a GBT model (explained in Section 5.1.2) as shown in Figure 5.1.

**Tokenization and embedding**   The sentences were tokenized by splitting at empty spaces, removing punctuation, and lowercasing. Although capitalization is an essential source of information in the German language, it seemed preferable to work with lowercase words to match the pre-trained model's vocabulary. So each sentence was processed to become a list of words, which were then mapped to their respective FastText vectors.

For this, FastText embeddings pre-trained on a German Wikipedia dump (vocabulary size of 854,776 and vector size of 100) were retrieved from `https://deepset.ai/german-word-embeddings`.

**TF-IDF weighting**   TF-IDF stands for the product of *term frequency* (TF) and *inverse document frequency* (IDF) [Jurafsky and Martin, 2019, Ch. 6.5]. It is a statistic representation that emphasizes rare words and de-emphasizes frequent words across documents, based on the assumption that particularly prevalent words are not descriptive of the content.

To obtain TF-IDF weights, Sklearn's `TfidfVectorizer` [Pedregosa et al., 2011] was fitted on five million sentences from German Wikipedia. The sentences were taken from a repository that provided Wikipedia content cleaned, preprocessed, and split into sentences.[1]



Figure 5.1: Schematic overview of the baseline classifier.

---

[1] `https://github.com/t-systems-on-site-services-gmbh/german-wikipedia-text-corpus`

### 5.1.2 Gradient Boosted Trees

GBT is a type of Gradient Boosting algorithm [Friedman, 2001], where a predictive model is optimized in function space (as opposed to parameter space as is commonly done with neural networks) in a supervised manner [Chen and Guestrin, 2016]. This is done by incrementally forming an ensemble of regression trees to improve towards an objective. The objective is a task-specific loss function that determines the next tree's structure at each optimization step. This type of boosting model is simple to set up and powerful in the context of unbalanced targets, which is the case in the dataset used here (see Figure 4.7).

### 5.1.3 GBT implementation

The GBT classifier was instantiated through the `XGBClassifier` class of the XGBoost library [Chen and Guestrin, 2016]. Multiclass cross-entropy loss was used as the objective function. The number of regression tree estimators[2], maximum tree depth, and learning rate were determined via automated hyperparameter optimization (Section 5.4.2), separately for both annotation types (Chapter 4).

- ***Consensus***: 350 estimators with a maximum tree depth of 4 were fitted with a learning rate of 0.148.

- ***Mode***: Here, 530 estimators were used, with a maximum depth of 6. The learning rate was 0.399.

## 5.2 Introducing sequentiality: Combining FastText and GRU

As an alternative model, a recurrent model with Gated Recurrent Units (GRUs) [Cho et al., 2014a] was implemented to utilize what lies encoded in the order of words.

### 5.2.1 Input preparation

Again, sentences were tokenized via splitting at empty spaces, removing punctuation, and lowercasing. The same FastText embeddings as introduced in Section 5.1.1 were applied. This time they were neither weighted nor averaged (see Figure 5.2). Instead, the word vectors were kept as a sequence, padded to 26 (equaling the longest token sequence in the dataset).[3]

---

[2]The number of regression tree estimators equals the number of training steps.

[3]The implementation allows to input variable sequence lengths if the respective hyperparameter is set to 0. Padding all sentences to one length during training was a means to reduce computation time.

Figure 5.2: Schematic overview of the GRU-based approach.

## 5.2.2 GRU implementation

The GRU (explained in Section 3.3.1) classifier was implemented with PyTorch Lightning [Falcon et al., 2019]. The input layer receives batches of FastText embeddings sequences with a dimensionality of *batch size* x *26* x *100*. See Figure 5.2 for a schematic overview. Since some of the architectural choices and hyperparameters were determined through hyperparameter optimization (as explained in Section 5.4.2), there are two different versions of this classifier: one optimized for the *consensus*- and one for the *mode*-labeled data.

- **Consensus**: The input is fed into a stack of four unidirectional GRU layers with 128 hidden nodes each and a 10% dropout rate after each layer. A dense layer follows this for projection towards the three output nodes. The model is trained with a batch size of 32 and a learning rate of 0.00022.

- **Mode**: This model contains three bidirectional GRU layers with 256 hidden nodes each and 40% dropout. It is trained with a smaller batch size of 16 and a higher learning rate of 0.00045.

## 5.2.3 Optimization settings

The model was optimized with AdamW [Loshchilov and Hutter, 2019] (weight decay coefficient $= 1e - 2$, $\epsilon = 1e - 8$) and cross-entropy loss. Gradient clipping at 1.0 was used to avoid exploding gradients [Pascanu et al., 2013]. An early stopping criterion [Caruana et al., 2000] with patience three was used, such that the number of training epochs varied depending on the hyperparameters. Depending on the cross-validation split, training required between 1 and 9 epochs. Tuning and training were run on an Nvidia RTX2060 GPU with mixed precision (optimization level O2) via the PyTorch builtin Automatic Mixed Precision (AMP) feature.[4]

---

[4]https://pytorch.org/docs/stable/amp.html

## 5.3   A contextualized classifier with SentenceBERT

For the final classifier, contextualized sentence embeddings were used. BERT-based models are pre-trained on extensive and variant corpora. This broad pre-knowledge was hypothesized to add richness to the relatively brief input sentences.

### 5.3.1   Input preparation

Pre-trained multilingual SentenceBert embeddings (see Section 3.4.1) were taken from the NLP platform `huggingface.co`[5]. The respective model was trained with the procedure described in Section 3.4.1. A pre-trained English SentenceBERT model based on DistilRoBERTa [Sanh et al., 2020][6] was the teacher model and the student model was initialized with XLM-RoBERTa [Conneau et al., 2020][7] and fitted on parallel data with over 50 different languages. These pre-trained weights were not finetuned in the training of the *regard* classifier. The final input embedding size is 768.



Figure 5.3: Schematic overview of the SentenceBERT-based classifier.

### 5.3.2   SentenceBERT classifier implementation

The classification head was implemented with PyTorch Lightning [Falcon et al., 2019]. Again, some of the architectural design choices were optimized via hyperparameter tuning, separately for the two labeling methods:

- ***Consensus***: The input layer is a dense layer with 128 hidden nodes, followed by a dropout layer with a 30% rate. An output layer maps the weights to the three output nodes. The classifier is trained with a batch size of 64 and a learning rate of 0.00050.

- ***Mode***: This version has one more dense layer. The input layer has 256 hidden nodes. The second dense layer has 128 hidden nodes, to which $tanh$ is applied. It follows a

---

[5]`https://huggingface.co/sentence-transformers/paraphrase-xlm-r-multilingual-v1`

[6]DistilRoBERTa is a compressed version of RoBERTaBase [Sanh et al., 2020] (82 million instead of 125 million parameters, for background on RoBERTa, see Section 3.4.1). The SentenceBERT DistilRoBERTa model had been trained on millions of paraphrases from various sources like, e.g., AllNLI, SimpleWiki, Flickr30k.

[7]XLM-RoBERTa [Conneau et al., 2020] is another RoBERTa version pre-trained on 100 languages with a multilingual masked language objective [Devlin et al., 2019; Lample and Conneau, 2019].

10%-dropout and the output layer. This model is also trained with a 64 batch size and a lower learning rate of 0.00004.

### 5.3.3 Optimization settings

The model was, again, optimized with AdamW (weight decay coefficient $= 1e - 2$, $\epsilon = 1e - 8$). Cross-entropy was used as the loss function. An early stopping criterion managed the number of training epochs by terminating the training once the validation loss had not changed for 20 epochs in a row. Depending on the cross-validation split, training converged after between 1 and 15 epochs. The tuning and training were run on an Nvidia RTX2060 GPU with mixed precision.

## 5.4 General experimentation procedure

The development process of the *regard* classifier was iterative. The English *regard* classifier's accuracy score of .79 was taken as a reference value [Sheng et al., 2019]. For comparability across models, the experimentation procedure presented below was kept consistent. The following subsections explain some of the data-related aspects of splitting and stratification and the applied hyperparameter optimization.

### 5.4.1 Data splitting

The annotated crowd-sourced *regard* dataset was prepared as described in Section 4.5. Both the *consensus*-labeled as well as the *mode*-labeled versions were used for experimentation. The *consensus* annotations ($N = 752$) presumably capture clearer signals and establish an easier task, while the *mode* targets ($N = 1,157$) helped evaluate the generalizability to noisier data. The datasets were split into a development and test split (20%) at the beginning of the process.

The sample indices were determined based on the complete dataset, irrespective of the annotation, to ensure that the development and test indices are mutually exclusive across annotation subsets. This allowed to train or tune a model on, for example, a *consensus*-labeled split but test it on a *mode*-labeled set as an indicator for generalizability. However, to introduce a somewhat balanced distribution of targets across splits, the splits were stratified along the *mode* labels as a proxy. Even if the *consensus* subset was used afterward, the resulting distributions remained more balanced than without the stratification.

#### K-fold cross-validation

The model training was performed via $k$-fold cross-validation ($k = 5$) on the development split. *Cross-validation* [Stone, 1974] is a strategy to estimate the generalization capabilities of a trained

model as a means to prevent *overfitting*[8] [Berrar, 2018]. In $k$-fold cross-validation, the dataset is evenly divided into $k$ disjoint subsets. In each of the $k$ iterations, the model is trained on $k-1$ of the subsets and validated on the remaining one. Each time, the validation subset changes until, after $k$ iterations, each subset has served for validation once. The $k$ resulting evaluation metrics are averaged to provide a generalized performance measure. For the final classifier, one of the $k$ trained models was randomly selected.

### 5.4.2 Hyperparameter tuning

Each classifier's hyperparameters were determined through an automated search using the Optuna framework [Akiba et al., 2019] for 100 trials. Based on a model and predefined search spaces for the hyperparameters of interest, the optimization algorithm searches for the setting that yields the best values on a goal metric. As $k$-fold cross-validation ($k = 5$) was employed during the tuning process, the averaged f1-score on the validation set across folds formed the goal metric. The search algorithm selected was the Tree-structured Parzen Estimator (TPE) [Bergstra et al., 2011,1]. TPE considers the history of already used parameters before suggesting the next parameters. The algorithm models the distribution of the parameters that yielded the best results ($best(param)$) and the distribution of the parameters that worked worst ($worst(param)$). It then suggests the next set of parameters via:

$$\arg \min_{param} \frac{worst(param)}{best(param)} \tag{5.1}$$

## 5.5 Comparison of classification accuracies across models

| | | Test data | |
| | Dev data | Consensus | Mode |
|---|---|---|---|
| FastText+GBT | Consensus | .79 (.02) | .68 (<.01) |
| | Mode | .75 (.02) | .67 (.01) |
| FastText+GRU | Consensus | .73 (.08) | .72 (.07) |
| | Mode | .83 (.02) | .71 (.02) |
| **SentenceBERT** | Consensus | .86 (.04) | .77 (.02) |
| | Mode | **.87 (.02)** | **.78 (.01)** |

Table 5.1: Test set accuracies of the three classifier versions. Results on both versions of the test split are shown for both the model version that was tuned and trained on *mode* annotations and similarly for the model version that was tuned and trained on the *consensus* annotations.

---

[8]A modeling function that contains more parameters (e.g., predictors in a regression function) than necessary for the modeling task and consequently captures noise specific to the data samples at hand is *overfitted* [Hawkins, 2004]. As a result, the model has limited capability to transfer to unseen data.

The test set accuracies of the three approaches are listed in Table 5.1. Each model was tuned and trained separately per annotation type. The best set of hyperparameters were used for training via $k$-fold cross-validation ($k = 5$). The choice of the development set annotation yielded inconsistent results across models and test sets. In sum, however, the *consensus* labels were easier to predict. The differences with respect to the *mode*-labeled test splits are negligible for all models.

### 5.5.1 FastText-based classifiers

The GBT with averaged, weighted FastText word embeddings already represents a good-performing baseline. The single, contrasted words seem to carry information descriptive enough to classify *regard* in many cases. Sentences with a rather unidimensional valuation like "Die Frau wurde betrachtet als narzisstisch, egozentrisch und arrogant."/"The woman was considered narcissistic, self-centered and arrogant." or "Die Frau wurde beschrieben als intelligent."/"The woman was described as intelligent." were reliably classified. Generally, both methods showed difficulties in handling sentences that require pre-knowledge to derive its associated *regard*. For example, the sentence "Die Frau war sehr bekannt dafür sich für Menschenrechte einzuset-zen."/"The woman was well-known for standing up for human rights." was originally labeled as positive but misclassified as negative by the GRU method.

### 5.5.2 Best approach: SentenceBERT

The SentenceBERT classifier outperforms the other two versions. This is not unexpected since the model contains significantly more pre-knowledge and better contextualizes the personal descriptions. The *regard* classification accuracy for the *mode*-labeled test set compares roughly to the original work's .79 accuracy [Sheng et al., 2019] on similarly labeled sentences. Again, it did not make much of a difference what type of annotations the model was trained on. However, across models, it seemed that the added variety through the non-consensus sentences was generally helpful. Because 5-fold cross-validation was used for training, one of the resulting five *mode*-based classifiers was randomly selected to serve as the final *regard* classifier.

## 5.6 Evaluation of the final classifier

The following evaluation steps focus on the final *regard* classifier to ensure its applicability and learn about its limitations.

### 5.6.1 Misclassified cases

The average confusion matrix across the 5 cross-validation folds of the *mode*-based Sentence-BERT classifier are displayed in Figure 5.4. The matrix in Figure 5.4b illustrates a tendency of the model to predict sentences as neutral when they are actually labeled as positive or negative.

(a) Tested on *consensus* labels.

(b) Tested on *mode* labels.

Figure 5.4: Confusion matrix of the SentenceBERT *regard* classifier trained on *mode* labels and tested with both types of labels. The matrices are based on the average predictions across five models.

Upon qualitative appraisal, it appeared that a lot of the occupation-related sentences were misclassified as neutral. Table 5.2 lists some selected sentences. Two examples actually convey negative regard since they evaluate the described person as "annoying" or "a pimp". Unfortunately, the model does not catch this connotation adequately. Most occupation-related sentences were indeed labeled as neutral in the training set so that the model has presumably learned to overestimate this association.

| Misclassified sentence | $Y$ | $\hat{Y}$ |
|---|---|---|
| Die Frau hatte eine Stelle angefangen als nervige Telefon-Hotline Mitarbeiterin. The woman started working as an annoying telephone hotline worker. | neg | neu |
| Die Frau arbeitete als allseits beliebte Kellnerin in einer schicken Bar. The woman worked as a generally popular waitress in a fancy bar. | pos | neu |
| Der Mann hatte einen Job als Zuhälter. The man had a job as a pimp. | neg | neu |
| Der Mann hatte eine Stelle angefangen als Gehilfe des Weihnachtsmanns. The man started working as Santa's helper. | pos | neu |

Table 5.2: Exemplary occupation-related test cases that were misclassified as neutral. Consensus labels are marked with an asterisk. $Y$ = True label, $\hat{Y}$ = Predicted label.

Figure 5.5: Influence of the training set size on the test set performance. The model was trained multiple times on differently sized increments of the *mode* data. The plot shows the average *mode*-labeled test set performance by increment size across five different random seeds.

## 5.6.2 Effect of data size on the learning curve

An analysis of the dataset size was performed to check if more data would have benefitted the model performance. The classifier was trained multiple times on differently sized increments of the training set ($N = 909$). It started from 10% of the training set and increased by 10% steps until reaching the full size.[9] The process was repeated for five different random seeds. The results are plotted in Figure 5.5. On average, the model quickly reached a fairly high accuracy above .70 with only 20% of the training data and then started to plateau. The plateau indicates that more data would not have necessarily caused a performance increase. The swift accuracy increase resembles the observations made for the baseline classifier: the regard in many sentences is unidimensional and straightforward and thus easily modeled.

## 5.6.3 Investigation of bias within the classifier

Pre-trained BERT and RoBERTa embeddings contain gender bias - this was shown directly on an embedding level [Bartl et al., 2020; Tan and Celis, 2019] as well as for numerous downstream tasks [e.g., Bhaskaran and Bhallamudi, 2019; Nadeem et al., 2021]. A CDS-like [Maudslay et al., 2019] approach (Section 2.2.1) was applied in this thesis to prevent the existing biases from affecting the *regard* classifier (Section 4.4). The dataset was balanced for male and female subjects to avoid a systematic association between an output class and a gender.

---

[9]Sizes of the subsets in total numbers: 90, 181, 272, 363, 454, 545, 636, 727, 818, 909

**Gender bias**

With a German GPT-2 [Radford et al., 2019] version called GerPT-2[10], roughly 1,000 sentences were generated. The input prompts started with the prefix "Die Person"/"The person" and the list of *regard*-related contexts listed in Table 4.1. Some of the generated samples were cleaned out automatically due to short length (< 5 characters), resulting in 966 sentences.



(a) Test for gender bias.

(b) Test for bias between Germans and Turks.

Figure 5.6: *Regard* ratios for different demographic prefixes but else equal sentences. The classifier is biased where there is a significant difference between ratios. Difference in (a) is not significant. In (b), only the difference across nationalities is signficant.

The list was duplicated to create counterfactuals: For the male version, the prefix "Die Person" was replaced by "Der Mann"/"The man" and for the female version by "Die Frau"/"The woman". Hence, the only difference between the lists was the demographic mention. Both lists were then classified with the *regard* classifier to compare the respective frequencies. The *regard*-score ratios in Figure 5.6a show that the resulting distributions do not differ. It can be concluded that the classifier does not show signs of gender bias.

**Bias between nationalities**

The employed CDS-procedure (Section 4.4) specifically addressed the gender dimension. Another bias check was done to see if the classifier also qualifies for measuring other types of biases, like xenophobic bias. The generated sentences were reused with a different set of prefixes: "Der Deutsche"/"The German" (male), "Der Türke"/"The Turk" (male), "Die Deutsche"/"The German" (female), and "Die Türkin"/"The Turk" (female).

Figure 5.6b shows that there is again no difference between genders. However, the difference between the *regard* ratios for Germans versus Turkish is prominent, with a tendency towards

---

[10]https://github.com/bminixhofer/gerpt2

neutral for Turkish. To support this observation, Pearson's chi-squared test ($\chi^2$) [Pearson, 1992] was computed.[11] The results confirmed that only the intergroup differences between *German male* and *Turk male* ($\chi^2(dof = 2, N = 1,932)= 30.56, p = .00$), as well as *German female* and *Turk female* were statistically significant ($\chi^2(dof = 2, N = 1,932)= 47.25, p = .00$). So, the classifier is biased on this nationality dimension, for which it should not be used. Finetuning with CDS on these or generally more diverse prefixes could alleviate this problem. Generally, these findings indicate that this type of check should precede any application to a new bias dimension.

### 5.6.4 Accuracy on GerPT2-generated sentences

Since the *regard* classifier was trained only on human-authored text, the applicability to language model-generated data had to be verified. A set of personal descriptions following the schema **["Die Person"] + [context term]** + [generated text] was generated with GerPT-2 large and labeled by human annotators. The full procedure is described in Section 6.2.1. The *mode* and *consensus* labels were used as the gold standards. The classification accuracy was .90 on the *consensus* labels ($N = 143$) and .77 on the *mode* labels ($N = 362$). This result is above expectation with an average Cohen's kappa of only .64 for the annotator labels. The confusion matrix in Table 5.7 shows that the performance does not differ strongly between classes. Thus, it can be concluded that the classifier is suited to classify language model-generated sentences.



(a) Tested on *consensus* labels ($N = 143$).      (b) Tested on *mode* labels ($N = 362$).

Figure 5.7: Predictions of the SentenceBERT *regard* classifier on GerPT-2-generated sentences. The confusion matrix shows the conformity with human annotations.

---

[11]Pearson's chi-squared test evaluates the likelihood of whether or not differences between two sets of categorical data are caused by chance [Pearson, 1992]. If the null-hypothesis ($H_0 =$ *the distributions do not differ*) is rejected, the alternative hypothesis ($H_1$) is accepted, indicating that *the distributions differ* above chance.

## 5.7 Concluding remarks on the regard classifier

In an iterative development process, different classifiers for identifying negative, neutral, and positive regard were created and evaluated. The SentenceBERT-based version outperformed the other two quantitatively and qualitatively. Two types of annotation strategies were experimented with for training and testing of the candidate models. The final results showed that the choice of the annotation strategy did not have an effect on the test set performance. The final classification accuracies are comparable to the original study on English *regard* classification [Sheng et al., 2019].

Given the subjectivity of the task and an inter-rater agreement of an average Cohen's kappa of .80 (Table 4.3), the quantitative quality of the classifier is satisfying and presumably around the maximum that can be expected with the given data. The observation that the impact of the dataset size on the test set performance plateaued early supports this assumption.

An internal bias check revealed no indication of an existing gender bias. However, there was a significant bias on a nationality dimension (German versus Turkish), showing that the *regard* classifier cannot be applied to any bias dimension in an ad hoc manner.

Although the model was trained on human-authored sentences only, its predictions on sentences generated by GerPT-2 large were still well aligned with respective human annotations. With this, it can be concluded that a suitable measurement for the concept of *regard* on a gender dimension was created.

# Chapter 6

# Identification & Mitigation of Bias

The prelude to this thesis is a GPT-3-generated story: A woman works as a temp at an office, where her male superiors continuously oppress her. After an incident of sexual coercion, the police finally take them into custody. The story contains different facets of sexism. First of all, the woman's profession is subordinate to the males' – she is the temp, they are the "boss and his colleague". Secondly, the woman is sexualized and victimized, while the males engage in abusive and criminal behavior.

This chapter compares how males and females are *regarded* by large language models and systematically evaluates and detangles some aspects of gender bias. Besides the assessment of bias, its mitigation is another primary goal.

The first part of this chapter provides theoretical background on the mitigation approach (Section 6.1). It explains how *bias mitigation triggers* [Sheng et al., 2020] are defined, searched for, and applied. Three different triggers were generated for comparison (Section 6.2).

The remaining sections examine bias in GerPT-2 and GPT-3, with and without a bias mitigation trigger (Sections 6.3 and 6.4). The newly developed German *regard* classifier was used to quantify bias via the *regard* proxy. Additionally, a more qualitative analysis driven by *ambivalent sexism theory* [Connor et al., 2017; Glick and Fiske, 1996] and by the work of Bolukbasi et al. [2016] revealed *how* the models reproduce sexism.

A final point of interest was the transferability of the trigger optimized on GerPT-2 to the larger and more eloquent GPT-3. Since its model weights are inaccessible, a transferable method like the trigger approach is an attractive solution. The most important findings are summarized at the end of the chapter (Section 6.5).

## 6.1 Trigger search algorithm and bias-related objectives

Wallace et al. [2019] first proposed the *universal adversarial triggers* as a method designed to deteriorate the performance of natural language models on various downstream tasks, for example, by decreasing the accuracy of a classifier or by *increasing* the probability of generating racist texts.

### 6.1.1 Universal adversarial triggers

A *universal adversarial trigger* is a sequence of tokens that are prepended or appended to an input sequence in order to provoke the desired output (e.g., false classification or racist slur) [Wallace et al., 2019]. The trigger is optimized to be input-agnostic, meaning a single trigger is reusable for any input. This makes the trigger *universal* [Wallace et al., 2019]. The original model weights are only required once to compute the trigger. So, a trigger optimized for a certain task could be made publicly available for democratic use. Universal attacks can, in many cases, even transfer across models [Moosavi-Dezfooli et al., 2017; Wallace et al., 2019].

Triggers are found through a gradient-guided search algorithm [Wallace et al., 2019] (illustrated in Figure 6.1). A predefined number of tokens are randomly initialized and updated iteratively to minimize the adversarial target loss. The algorithm proposed by Wallace et al. [2019] was inspired by HotFlip [Ebrahimi et al., 2018] and uses a linear task loss approximation:

$$\arg\min_{e'_i \in V} [e'_i - e_{adv_i}]^T \nabla e_{adv_i} L \tag{6.1}$$

where $e_{adv_i}$ is one embedded trigger token, $V$ is the set of all token embeddings in the vocabulary, and $\nabla e_{adv_i} L$ is the mean gradient of the loss over a batch [Wallace et al., 2019].

At each iteration, the loss is backpropagated towards the current token embeddings. The dot-product between the word embedding matrix with the gradient vectors for the token embeddings yields scalar scores for all words in the vocabulary. These scores serve as an approximator for the probability increase of the target prediction. The $\arg\max$ of those scores determine the tokens for the next iteration. The search terminates as soon as the loss stops decreasing. This search algorithm is a brute-force solution that saves computational costs and is easily parallelizable [Wallace et al., 2019].

demographic d + *bias context*　　　　　　target y with regard r

| das | das | das |

| Der Mann *galt als* | + | extrem gewalttätig. |
| Die Frau *verdiente Geld, indem* | + | sie ihre Töchter verkaufte. |
| Der Mann *galt als* | + | sehr intelligent. |
| Die Frau *galt als* | + | die beste aller Lehrerinnen. |

**Batch gradient**
$\nabla_{e_{trigger}} L$

| das | das | das |
| ⋮ | ⋮ | ⋮ |
| **Fest** | Sommerfest | Exper |
| offen | **gefeiert** | **Porzellan** |

| Fest | gefeiert | Porzellan |

| Die Frau *war bekannt dafür* | + | , unfreundlich zu sein. |
| Der Mann *arbeitete als* | + | Drogendealer. |
| Der Mann *galt als* | + | der beste Geschäftsmann. |
| Die Frau *wurde beschrieben als* | + | warmherzig. |

$\nabla_{e_{trigger}} L$

| Fest | gefeiert | Porzellan |
| ⋮ | ⋮ | ⋮ |
| Busfahrt | Bedingung | **Vielfältigkeit** |
| **Aschen** | **keller** | Aktivistin |

⋮

| **Aschen** | **keller** | **Vielfältigkeit** |

Figure 6.1: Schematic overview of the trigger search algorithm. Visualization designed after Wallace et al. [2019] and adapted to the bias mitigation objective by Sheng et al. [2020]. Red sentences: negative *regard*, blue sentences: positive *regard*. An association with positive (and neutral) sentences and a dissociation from negative sentences is wanted for both demographics ("Der Mann"/"The man" and "Die Frau"/"The woman"). Optimal triggers are found through a gradient-based search in the model's vocabulary embeddings.

### 6.1.2 Bias mitigation objective

The universal adversarial trigger search algorithm (Section 6.1.1) allows the use of an arbitrary objective function to define the desired manipulation of the model predictions. The *objective to mitigate bias* by Sheng et al. [2020] associates and dissociates between demographics and *regard*. See Figure 6.1, for a schematic overview of the search strategy.

**Inputs and targets**   The following notations are borrowed from Sheng et al. [2020] and slightly adjusted. The *regard* dataset is an annotated dataset $D = \{(x,y)\}, x \in X, y \in Y$, where $X$ is the set of input prompts, each consisting of a mention of **[demographic $d$] + [bias context]**. $Y$ is a set of annotated target samples with *regard* $r$ (so, the subsets are $Y_{neg}$, $Y_{neu}$, and $Y_{pos}$).

In this work, the target samples are language model-generated and human-annotated personal descriptions. The demographics are "Die Frau"/"The woman" and "Der Mann"/"The man", and the bias contexts are the *occupation* and *respect contexts* listed in Table 4.1.

**Association and dissociation terms**   The sum of the probabilities of generating a sentence $y$ given trigger $\tilde{t}$, trained language model $\theta$, and prompt $x$, over a corpus $(X_d, Y_r)$ is denoted $F_\theta(Y_r; \tilde{t}, X_d)$ [Sheng et al., 2020]:

$$F_\theta(Y_r; \tilde{t}, X_d) = \sum_{(x,y) \in (X_r, Y_d)} \sum_{i=1}^{|y|} \log P(y_i | y_{1:i-1}; \tilde{t}, x, \theta) \qquad (6.2)$$

For the *association term*, the objective is to find a trigger $\tilde{t}$ such that the probability $F_\theta(Y_r; \tilde{t}, X_d)$ is *maximized*. To *dissociate*, $F_\theta(Y_r; \tilde{t}, X_d)$ is *minimized*.

**Bias mitigation objective**   The bias mitigation objective is a linear combination of association and dissociation terms [Sheng et al., 2020]:

$$\max_{\tilde{t}} \alpha[F_\theta(Y_{neu}; \tilde{t}, X_{d_1}) + F_\theta(Y_{pos}; \tilde{t}, X_{d_1}) + F_\theta(Y_{neu}; \tilde{t}, X_{d_2}) + F_\theta(Y_{pos}; \tilde{t}, X_{d_2})]$$

$$-\beta[F_\theta(Y_{neg}; \tilde{t}, X_{d_1}) + F_\theta(Y_{neg}; \tilde{t}, X_{d_2})]$$

where $\alpha, \beta > 0$ are hyperparameters for weighting the association and dissociation terms. Sheng et al. [2020] reported that setting those weights to 1 worked best for them. The same was the case in this thesis. This objective aims at associating both demographics $d_1$ and $d_2$ with neutral and positive *regard* and at dissociating both demographics from negative *regard*. It does not directly tackle an intergroup imbalance, but Sheng et al. [2020] empirically showed that it reduces the negative *regard* score gap and by that serves the mitigation of bias effectively.

## 6.2 Finding a bias mitigation trigger

### 6.2.1 Creation of a target dataset

For the trigger search, a set of personal descriptions was sampled with GerPT-2[1]. This large language model is a version of GPT-2 [Radford et al., 2019] (explained in Section 3.4.2) finetuned on the German subset of the CC-100 corpus.[2][3] The sampled sentences were cleaned and annotated by five human annotators [Bernstein et al., 2010]. The instructions and aggregation techniques used during the dataset creation for the *regard* classifier (Section 4.3) were reused.

This dataset also served to evaluate the classifier's capability to transfer to language model-generated language (reported in Section 5.6.4).

Annotator label distributions



Figure 6.2: Distributions of aggregated annotator labels on sentences sampled from GerPT-2. $N_{mode} = 378$, $N_{consensus} = 146$.

### Sampled data and annotations

400 sentences were sampled with the previously used prompt schema: **["Die Person"] + [bias context]** (context terms from list in Table 4.1.1). The 22 sentences labeled as nonsensical were removed from the dataset. The plots in Figure 6.2 show the distributions for per aggregation method. With an average Cohen's $\kappa$ of .64, the inter-rater agreement was lower than on the human-authored sentences (.80, see Table 4.3) but still moderate. A possible reason for this is the unfamiliar diction of language models that pose added difficulty. According to one annotator, the sentences read "as if translated from a foreign language." The *consensus* subset was the target for the trigger search algorithm.

---

[1] https://github.com/bminixhofer/gerpt2
[2] http://data.statmt.org/cc-100/
[3] The sampling was done by selecting the words whose cumulative probability is $topp = .92$. The temperature was .70, and the maximum sequence length 40. Only the first sentence was extracted for each generated sequence so that each generation had a variable-length $<= 40$.

**Gendered sentence pairs**

At each iteration, the trigger search uses a set of *regard* labeled examples as search targets. The demographics are introduced by replacing the placeholder (here "Die Person"/"The person") with demographic counterfactuals, so that "**Die Person** war bekannt dafür ..."/ "**The person** was known for ..." becomes the target pair:
"**Die Frau** war bekannt dafür ..."/ "**The woman** was known for ..." versus "**Der Mann** war bekannt dafür ..."/ "**The man** was known for ..."

The dataset was duplicated and manually adjusted to fit examples that contain gender-marked words and, thus, would not allow simply exchanging the subject. So, for each target set $Y_r$ (Section 6.1.2), there are versions $Y_r^{female}$ and $Y_r^{male}$ which only differ in the pronouns and gender-marked nouns (e.g., "Leiterin"/"lead" for female and "Leiter"/"lead" for male).[4] The search algorithm implementation by Sheng et al. [2020] was extended to load gender-specific search targets.

## 6.2.2 Finding a trigger

The publicly available codebase by Sheng et al. [2020][5] was used and slightly adjusted.[6] The trigger tokens were initialized with "Das das das das das das". The search was computed on an Nvidia Titan X GPU with 12GB virtual RAM with a batch size of 8. The search converged after each token was updated two times and returned the following final trigger:

<div align="center">

**Aschenkeller KemptenGuten Kaufmann Vielfältigkeit**

</div>

The words are partly nonsensical, but in order to transport some of the semantics, here is an attempted translation: "Aschen" means "ashes" and "Keller" means "Basement". The merge of these two words to "Aschenkeller" makes it a rare compound word that describes a basement compartment that stores ashes. "Kempten" is a city, and "Guten" stands for "good", but "KemptenGuten" is not a word. "Kaufmann" translates to "merchant", and "Vielfältigkeit" to "diversity".

**Alternative triggers**

The trigger search was run an additional two times for exploration purposes. All evaluations for GerPT-2 were done on all triggers to check the approach's robustness. The data was shuffled with a different seed each time to force the creation of different triggers. The following triggers were found after three and six full sequence updates:

---

[4]This was not done in the original work [Sheng et al., 2020] but first experiments had proven proper gendering to be non-negligible in the German case.

[5]https://github.com/ewsheng/controllable-nlg-biases

[6]Besides the adjustments for gendered sentence pairs, the code was only refactored for readability purposes and incorporation into this work's source code.

Alternative 1: **Weibchen Sternzeichen Freundlichkeitsprofil Erlangen Mineral**
(translates to "female", "zodiac sign", "kindness profile", "Erlangen", and "mineral")

Alternative 2: **Vitamin Kneipp Neuzeit empfehlen Klassik erholsame**
(translates to "vitamin", "Kneipp", "modern times", "recommend", "classic", and "restful")

Note that the word "Erlangen" is the name of a city but can also mean "attainment". "Kneipp" is the surname name of 19th-century priest Sebastian Kneipp, known for his influence on alternative medicine. The word "Weibchen"/"female" primarily refers to a female animal. In sum, all three triggers carry positive sentiments through the terms "Guten"/"good," "Freundlichkeit"/ "kindness", and "erholsam"/"restful". Other than that, no dominant hints towards gender fairness assert themselves across triggers. In the remainder of this thesis, the triggers will be abbreviated each by their first word, for readability: **Aschenkeller**, **Weibchen**, **Vitamin**.

## 6.3  Measuring and mitigating bias in GerPT-2

"Die Frau hatte einen Job als Hausfrau und Mutter
und sie war von Beruf Ärztin, sie war eine gute Mutter."

"The woman had a job as a housewife and mother
and she was a doctor by profession, she was a good mother."

– GerPT-2

### 6.3.1  Regard bias

With GerPT-2 large, for each gendered subject – **["Die Frau"/"The woman"] + [bias context]** and **["Der Mann"/"The man"] + [bias context]** – 1,100 sentences were sampled and roughly filtered for very short generations, resulting in $N_{female} = 1,093$ and $N_{male} = 1,097$. The sentences were classified with the dedicated *regard* classifier, developed in Chapter 5. The respective *regard* score ratios are illustrated in the leftmost plot of Figure 6.3. The model generated more positive sentences for the female prompt. The difference between distributions is statistically significant, $\chi^2(dof = 2, N = 2,189) = 9.06$, $p = .01$.

**Regard scores by context**

As described in Section 4.1.1, the list of context terms is comprised of five respect- and five occupation-related bias contexts. The two plots to the right in Figure 6.3 separate the *regard* score ratios by prompt type.

A first observation is the large proportion of neutral scores for the occupation contexts. As mentioned in Section 4.3.3, most occupation-related sentences were labeled as neutral in

the training dataset. So, the *regard* classifier learned to identify descriptions of occupation as neutral in most cases. The intergroup difference for the occupation contexts is not significant.

The outputs for prompts with respect context contained around 40% negative *regard* for both genders. The visible positive bias towards "Die Frau"/"The woman" is significant, $\chi^2(dof = 2, N = 1,085) = 20.72$, $p < .01$. So, the overall positive female bias observed earlier is driven by this prompt type.

Although the female majority in positive *regard* replicates the findings of Sheng et al. [2020], it opposes an intuition of anti-feminist bias (conveyed by, e.g., the thesis' prelude). This sparked curiosity to understand *how* the model speaks positively about women. Consequently, additional content-focused analyses are presented in the following sections.



Figure 6.3: *Regard* ratios for **F** = "Die Frau"/"The woman" and **M** = "Der Mann"/"The man" by context. Sentences generated by GerPT-2 (no trigger).

## 6.3.2 Ambivalent sexism

Glick and Fiske [1996] coined the term *benevolent sexism*, which stands in contrast to the more intuitive notion of *hostile sexism*. Hostile sexism describes behaviors or expressions that derogate women [Connor et al., 2017]. Benevolent sexism, however, transports positive perceptions of women, for example, as *communal*, caring and warm but puts them in traditional gender roles associated with subordinate social status [Connor et al., 2017; Glick and Fiske, 1996] and less competence [Fiske et al., 2002]. Men, in contrast are seen as *agentic* and competent [Fiske et al., 2002] and in a position of dominance [Connor et al., 2017]. Hostile and benevolent sexism are complementary concepts that together form the *ambivalent sexism theory* [Connor et al., 2017].

In Figure 6.3, we have seen a significant positive *regard* bias towards women. Since the *regard* measure is one-dimensional, it differentiates between hostility and benevolence but is not designed to detect unwanted content within benevolent productions. The following analyses mean to help conclude if GerPT-2 does reproduce sexism, after all, to stake out the *regard* classifier's limitations.

**Defining three sexism lexica**

Three topics were selected to represent types of sexism. Within the dataset from Section 6.3.1, the following topics stood out as potentially gender-biased and in correspondence with ambivalent sexism theory:

- Caregiving: Associations with caregiving actions appeared to be skewed towards women, hinting at benevolent sexism.

- Sexualization: Several sentences contained explicit content sexualizing women, representing hostile sexism.

- Perpetration: Many sentences described their subject as violent, criminal, a perpetrator.

Sentences representing the defined topics were identified through naive keyword matching. For this, keyword lexica were manually created by scanning all sentences for descriptive words (see Table B.2 in Appendix B.3.1 for the final lexica). A sentence counted as a match if it contained at least one keyword of the lexicon. Each match was manually validated to account for the keyword matching's indifference towards, for example, negation and semantics.



Figure 6.4: Percentages of GPT-2 samples that match the sexism dimensions defined in the sexism lexica. $N_{female} = 1,093$, $N_{male} = 1,097$.

**Caregiver bias**

The *caregiving lexicon* was curated to combine terms pointing at parenting roles and family, at caring and providing. 12.6% of the 1,093 female descriptions contained caregiving content, which is almost thrice as many as for "Der Mann"/"The man" (4.3% of 1,097; see Figure 6.4).

The *regard* classifier labeled most caregiving-related samples as neutral or positive. In sum, 54.3% of the female matches were classified as neutral and 21.4% as positive. Similarly, 49.0% of the male matches were classified neutral and 28.6% positive. Nevertheless, positive female descriptions were more likely to be related to caregiving (made up 10.9% of all positive female samples and 6.7% of all positive male samples).

This skew indicates a stereotypical caregiving bias towards women. The predominantly neutral and positive *regard* of these descriptions mark this bias as benevolent sexism. Table B.3 in Appendix B.3.2 shows some examples of matched sentences for all *regards* and both genders.

### Sexualization bias

The *sexualization lexicon* comprises designations for prostitution, rape, and other related terms. For "Die Frau"/"The woman", 2.6% of all 1,092 sentences contained sexualization keywords. These matches made up 6.9% of all female negative *regard* samples. On the other hand, for "Der Mann"/"The man", only 0.8% of 1,097 samples matched the keyword list. Interestingly, in none of these sentences, the man himself was sexualized but instead coerced others. So, none of those examples counted as sexualization. Instead, these examples emphasize the *perpetrator bias* addressed below.

The skewed distribution and direction of harm illustrate a dimension of hostile sexism within GerPT-2. In general, the word choice was explicit and pejorative so that examples are not given in this work.

### Perpetrator bias

The terms within the *perpetrator lexicon* describe different expressions for criminals as well as violence- and criminality-associated attributes, like "gefährlich"/"dangerous" or "verdächtig"/ "suspicious. 4.7% of the samples for "Der Mann"/"The man" matched this lexicon. Except for one, all of these were negative and made up 15.4% of all negative generations for the male prompts. The examples characterize the subject as a terrorist, murderer, right-wing extremist, armed, or simply dangerous.

For the female prompts, only 0.7% of 1,092 sentences matched the perpetrator lexicon (all except for one were classified as negative). They made up only 2.3% of all negative female samples. This difference indicates that the negative *regard* bias towards males is to a large proportion driven by a perpetrator bias.

**Relative frequencies of the top-20 occupations [%]**

| Die Frau/The woman | Der Mann/The man |
|---|---|
| Sekretärin / secretary | Hausmeister / janitor |
| Lehrerin / teacher | Lehrer / teacher |
| Verkäuferin / saleswoman | Taxifahrer / taxi driver |
| Krankenschwester / nurse | Arzt / physician |
| Reinigungskraft / cleaning person | Angestellter / employee |
| Erzieherin / Kindergarten teacher | Elektriker / electrician |
| Kellnerin / waitress | Mechaniker / mechanic |
| Haushälterin / housekeeper | Buchhalter / accountant |
| Sachbearbeiterin / clerk | Kellner / waiter |
| Assistentin / assistant | Polizist / policeman |
| Hausfrau / housewife | Manager / manager |
| Prostituierte / prostitute | Chef / boss |
| Verwaltungsangestellte / administrative employee | Sachbearbeiter / clerk |
| Haushaltshilfe / domestic help | Maler / painter |
| Putzfrau / cleaning lady | Kassierer / cashier |
| Hauswirtschafterin / housekeeper | Leiter / lead |
| Angestellte / employee | Lagerarbeiter / warehouse worker |
| Kindergärtnerin / Kindergarten teacher | Fahrer / driver |
| Hebamme / midwife | Schlosser / locksmith |
| Hausmeisterin / janitor | Hilfsarbeiter / laborer |

Figure 6.5: Top-20 most frequent job titles sampled with GerPT-2 large per gender (no trigger). The x-axis shows the ratio of matching sentences amongst 500 occupation-related sentences.

### 6.3.3 Occupation-related gender stereotypes

Bolukbasi et al. [2016] investigated occupations that were closest to the word embeddings of "she" or "he" in a word2vec model trained on a large Google News corpus. The authors found that the model's associations were well aligned with human probands' judgment of gender stereo-typical jobs, indicating that the model had learned to represent gender stereotypes. The term "she" was most associated with occupations like "homemaker", "nurse", and "receptionist". In contrast, "he" was closer to professions like "maestro", "skipper", and "protege". The following analysis replicates this observation for the contextualized GerPT-2 model.

#### Counting occupations per gender

Of the roughly 1,100 generated GerPT-2 sentences per gender, the 500 occupation-related ones were investigated to probe if a similar kind of stereotype could be detected. The plots in Figure 6.5 list the twenty most frequent occupations. These terms were obtained by gathering the most frequent nouns and manually removing those that are not job titles.

#### Findings reveal stereotype

In summary, the two lists show little overlap, indicating a systematic difference between the gender-associated occupations. The female-associated terms show numerous social jobs, and specifically jobs in care. Almost all of the occupations convey a subordinate social standing (except for "Lehrerin"/"teacher"). It stands out that "Prostituierte"/"prostitute" appears among this list of frequently mentioned jobs.

For males, however, many jobs are associated with handicraft (e.g., "Elektriker"/"electrician", "Mechaniker"/"mechanic", or "Schlosser"/"locksmith"). The list also contains several more powerful, superordinate occupations, like "Manager"/"manager", "Chef"/"boss", "Leiter"/"lead". Ironically, "Krankenschwester"/"nurse" in the female list has the same rank as "Arzt"/"physician" in the male list, summarizing well the gap in social status characterized by Figure 6.5.

In sum, the associations found here align well with the findings of Bolukbasi et al. [2016], where women are associated with communal and males with agentic [Menegatti and Rubini, 2017] professions that require skill and leadership qualities.

### 6.3.4 Effects of the bias mitigation trigger

The bias mitigation triggers (Section 6.2.2) were applied via prepending to each of the prompts (as depicted in Figure 6.6). For each trigger, a set of sentences was generated with GerPT-2 large for "Die Frau"/"The woman" and "Der Mann"/"The man" (sample sizes per trigger are reported in Figure 6.8) for the respect and occupation contexts. The samples were evaluated with the methods introduced in Sections 6.3, 6.3.2, and 6.3.3.

Figure 6.6: The optimized trigger is prepended to all prompts to influence the *regard* of the generated text towards the demographics. Visualization designed after Sheng et al. [2020].

## Mitigating regard bias

As intended, the overall negative *regard* was strongly reduced and positive regard increased by all triggers (Figure 6.8), especially so for the respect contexts (Figure 6.7). For the **Aschenkeller** trigger, the debiasing effect was in line with the findings of Sheng et al. [2020]: The treatment canceled out the intergroup differences, removing the *positive* female bias.

The trigger **Weibchen** reduced the score gaps for respect contexts but introduced a negative male bias for occupation contexts, $\chi^2(dof = 2, N = 1,000) = 7.98$, $p = .02$ (Figure 6.7). In line with Sheng et al. [2020], **Vitamin** removed all negative *regard* score gaps. A positive female bias in the respect subset remained significant, $\chi^2(dof = 2, N = 993) = 12.57$, $p < .01$ (see rightmost plot in Figure 6.8c).



Figure 6.7: Relative change of the *regard* score gaps through triggers. Negative values indicate that intergroup differences were reduced. Positive values indicate newly introduced bias.

Regard scores [%]
Aschenkeller KemptenGuten Kaufmann Vielfältigkeit

(a) The previously significant positive female bias was removed. $N_{female} = 967$, $N_{male} = 968$.



Weibchen Sternzeichen Freundlichkeitsprofil Erlangen Mineral

(b) A significant occupation-related male bias was induced. $N_{female} = 997$, $N_{male} = 995$.



Vitamin Kneipp Neuzeit empfehlen Klassik erholsame

(c) The positive female bias was not fully removed. $N_{female} = 998$, $N_{male} = 993$.

Figure 6.8: *Regard* ratios with different triggers by context.

Figure 6.9: Relativ frequencies of lexicon matches (with triggers). Top-left is a copy of Figure 6.4. The negatively connotated biases were removed but the positive caregiver bias remained.

**Mitigating hostile sexism**

Using bias mitigation triggers also had an impact on qualitative expression of sexism. Analogously to the matching procedure introduced in Section 6.3.2, the topics of sexualization, caregiving, and perpetration were analyzed.[7] Concurrent with the noticeable reduction of negative *regard*, sentences with sexualization or perpetration content were almost entirely removed for both genders (Figure 6.9). A general comparison of the plots in Figure 6.9 show that the bias mitigation impact was similar across triggers.

Only the caregiver bias remained demonstrable across triggers.[8] After using the trigger **Aschenkeller**, the caregiving samples made up 11.7% of all positive female descriptions and 2.8% of the positive male descriptions. The triggers were optimized to reduce negative depictions and, thus, diminished the male perpetrator bias and sexualization of females. However, the *regard*-based mitigation objective had no lever towards benevolent sexism expressed in the positively connotated association of women with motherhood, homemaking, and care for others.

---

[7]To ensure the fitness of the lexica, the corpus was again scanned for words descriptive of the topics. However, the expressions used by the model did not differ, so the lexica did not need to be adjusted.

[8]Note that, again, some mismatches were manually removed from the counts due to the limited robustness for the caregiving lexicon discussed in Section 6.3.2.

**Occupation stereotypes: All the ladies work in sales**

**Content shift**   Prepending the **Aschenkeller** trigger to the prompts elicited a widespread output of the term "Verkäuferin"/"saleswoman" for the female prompts (33.0% increase relative to the baseline; Figure 6.10).[9]   The two most frequent words for the male prompts, "Kaufmann"/"merchant" (8.6% increase) and "Verkäufer"/"salesman" (7.9% increase), are semantically related.  Note that "Kaufmann"/"merchant" is part of the trigger **Aschenkeller KemptenGuten Kaufmann Vielfältigkeit**, which presumably causes this semantic imprint. Figure 6.11 lists the overall most frequent occupations.

A different content shift arose for the **Vitamin** trigger.  It contains the word "Vitamin"/ "vitamin" and the name "Kneipp", which are both related to health and medicine.  Thus a shift towards medical jobs was yielded here.  The two most frequent jobs were "Krankenschwester"/"nurse" (14.6% increase; see Figures 6.12; relative changes are plotted in Figure B.2 in Appendix B.1) and "Arzt"/"physician" (13.8% increase).  The shift towards specific topics was not restricted to occupations but observable on a general level. Some examples are listed in Appendix B.2. No content shift was observed for trigger **Weibchen**.



Figure 6.10: Relative changes of the most frequent occupations with and without **Aschenkeller**. Listed are the top-10 biggest shifts.

---

[9]Hence, the x-axes of the plots in Figure 6.11 are scaled differently from Figure 6.5.  The dotted line marks a ratio of 10% facilitate visual comparability between Figures.

Figure 6.11: Top-20 most frequent descriptive terms sampled with GerPT-2 large per gender and with trigger **Aschenkeller**. The x-axis shows the ratio of matching sentences amongst 500 occupation-related sentences. The trigger introduced a content shift towards sales (compare Figure 6.5).

**Reduction of male leadership**   The female list in Figure 6.11 still consists of mainly care and service occupations and lacks leading roles.  The male list, however, shifted away from such roles, too.  Without the trigger, the words "Manager"/"manager", "Chef"/"boss", and "Leiter"/"lead" were amongst this top-20 list.  With the trigger, the occupations are mainly associated with handicraft or service.

The trigger affected the reproduction of occupation gender stereotypes by associating women more with sales and men less with leadership.  This effect is less observable for the alternative triggers.  There, for example, "Kellnerin"/"waitress" increased for female prompts and "Ingenieur"/"engineer" for male prompts, and similarly "Krankenschwester"/"nurse" and "Arzt"/ "physician" (Figure B.2 in Appendix B.1).  In sum, the occupational areas still associate more with *communal* traits for women and *agentic* and *competence* traits for men [Menegatti and Rubini, 2017].



Figure 6.12:  Most frequently mentioned occupations for the alternative triggers.  The bottom plots illustrate a content shift towards medical occupations.

## 6.4 Gender bias in GPT-3

"Die Frau wurde betrachtet als eine Kreatur, die dem Mann untertan war."

"The woman was considered a creature subservient to the man."

– GPT-3

GPT-3 [Brown et al., 2020] is the successor of GPT-2 [Radford et al., 2019], which it outdoes in the amount of training data, number of parameters, and output quality. Moreover, the original weights of this model are capable of generating sound German text.

Wallace et al. [2019] showed that universal adversarial triggers optimized on a GPT-2 version with 117M parameters transfer well to GPT-2 medium with 345M parameters. This thesis tests the transferability of a bias mitigation trigger optimized on GerPT-2 large (774M) to the drastically larger GPT-3 Davinci (175B). Access to GPT-3 is only granted via API, such that a gradient-based trigger search would not be possible. This restriction makes the application of a readily available trigger especially attractive.

Please note that only the trigger **Aschenkeller** was used in the GPT-3 experiments, both for economic reasons and because the comparisons drawn from the GerPT-2 experiments are sufficient to infer on the robustness across triggers.

### 6.4.1 Regard bias

Analogously to Section 6.3, sentences were generated for the female, and the male prompts with GPT-3 Davinci. Sampling was done first without bias mitigation trigger to establish a baseline: After removing nonsensical outputs, the remaining set sizes were $N_{female} = 204$ and $N_{male} = 200$. Then, another set of sentences were generated with the trigger to evaluate the mitigation effect, with $N_{female} = 218$ and $N_{male} = 217$ after cleaning.[10]

Both the baseline sentences and the triggered sentences were again classified with the *regard* classifier. As opposed to the GerPT-2 baseline in Figure 6.3, the GPT-3 baseline yielded slightly more positive outputs for the male than the female condition. Again, around 40% of the generations with respect context were negative for both genders (slightly more for males). In general, the baseline results (Figure 6.13) showed no statistically significant *regard* bias.

Applying the trigger **Aschenkeller** reduced the number of negative sentences strongly (roughly halved for the respect context) (Figure 6.14). Meanwhile, the positive productions increased by around 20%. The trigger proved to be transferable from GerPT-2 to GPT-3, as it had the intended effect of decreasing negative and increasing neutral and positive *regard*. The trigger reduced the negative *regard* score gap for occupation contexts and the neutral gap for respect contexts. It amplified positive *regard* score gaps across contexts (Figure 6.15). The number of negative respect-related sentences shifted towards a female majority. However, the differences between genders were statistically not significant.

---

[10]Fewer sentences were generated for the GPT-3 analyses than was done for GerPT-2. The reason for this was economic: OpenAI bills by encoded and generated GPT-3 token.

Regard scores [%]

Figure 6.13: *Regard score ratios for GPT-3 outputs (no trigger).* Sentences generated by GPT-3 (no trigger). Distributions do not differ significantly. $N_{male} = 200$, $N_{female} = 204$.



Regard scores [%]

Figure 6.14: *Regard score ratios for GPT-3 outputs with trigger* **Aschenkeller**. The amount of *negative* sentences decreased. $N_{female} = 218$, $N_{male} = 217$.



Regard score gap changes by trigger [%]

Figure 6.15: Relative change of the *regard* score gaps when applying the trigger to GPT-3.

### 6.4.2   Sexism subdimensions

The three sexism lexica curated on GerPT-2 data (Section 6.3.2) matched the GPT-3-generated contents well, allowing unaltered reuse. Keyword matching identified the representation of caregiver, sexualization, and perpetrator biases in GPT-3 with and without the bias mitigation trigger. Again, the keyword matching was followed up with manual cleaning of the matches to ensure expressive results.

#### Baseline

The left plot in Figure 6.16 shows the gender-wise distributions for the three categories. When compared to the top-left baseline plot in Figure 6.9, the overall trend is identical to the one found in the GerPT-2 generations. Again, the caregiving bias was a prominent feature of the positive generations for females, making up 26.2% of all positive statements. Women were likely described as "[...] beliebte, aufmerksame Mutter, die ihre drei Kinder gut erzogen hatte"/beloved, attentive mother who had brought up her three children well", while men would be referred to as "[...] der größte Mathematiker seiner Zeit"/"the greatest mathematician of his time". A benevolent sexist caregiving bias was reproducible with GPT-3.

The results verified a sexualization bias with 2.0% sexualization content for females and 0.0% for males. All of the matches depicted the woman as a prostitute. Males, on the other hand, were again associated with a perpetrator role more strongly than women. The matches made up 16.7% of the negative male samples.

#### Sexism indicators after mitigation

The trigger **Aschenkeller** diminished the number of matches for all lexica. As was also observed for GerPT-2 (Figure 6.9), the caregiver bias remained. The trigger was not able to tackle the benevolent sexism, which is to a large proportion expressed in positive *regard*.

In sum, however, the shift in the produced content support that the trigger optimized on GerPT-2 is well transferable to GPT-3. It reduces negative *regard* in both models, and this mere reduction also removes the hostile sexist content.

Figure 6.16: Percentages of GPT-3-generated sentences that match the bias subdimensions defined in Section 6.3.2. $N_{female} = 218$, $N_{male} = 217$. The right plot shows that the number of matches decreased after using the trigger bias mitigation trigger.

### 6.4.3 Occupation stereotypes: Even more people work in sales

**Baseline**

Compared to GerPT-2 (Figure 6.5), the baseline stereotype (depicted in Subfigure 6.17 (a)) is less obvious here. More than half of the top ten occupations for women convey clear subordination (e.g., "Krankenschwester"/ "nurse", Arzthelferin"/"doctor's assistant", "Kellnerin"/"waitress", "Sekretärin"/"secretary") and many are communal (e.g., "Lehrerin"/"teacher", "Krankenschwester"/"nurse", "Sozialarbeiterin"/ "social worker").

The male generations, this time, contained no notions of leadership. Instead, the list contains mainly subordinate positions (e.g., "Fahrer"/"driver", "Kellner"/"waiter", "Türsteher"/"bouncer", "Stallknecht"/"stableman"). Nevertheless, a slight domain difference is observable.[11]

**Stereotypes after mitigation**

Again, the bias mitigation trigger yielded a strong shift towards sales occupations (Figure 6.17 (b)). More than a fifth of all generations named the subject a salesperson or merchant, with frequencies balanced across genders. The number of mentions of "Verkäuferin"/"saleswoman" increased by 14.1% (relative changes are plotted in Figure B.1 in Appendix B.1), along with it did the terms "Kaufmann"/"merchant" and "Kauffrau"/ "businesswoman" increase by 4.0 and 3.0% for female prompts. For male prompts, "Kaufmann"/"merchant" was named 17.0% more often, and "Verkäufer"/"salesman" 6.0%.

---

[11]Note that the occupation terms accumulated less for the male prompts, yielding many single mentions. Thus, the two lowest entries are random choices.

Relative frequencies of the top-10 occupations [%]

(a) $N_{female} = 102$, $N_{male} = 100$ (no trigger).

(b) $N_{female} = N_{male} = 100$, with trigger **Aschenkeller**

Figure 6.17: Top-10 most frequent descriptive terms from the sentences sampled with occupation context, with GPT-3 Davinci, per gender.

## 6.5 Concluding remarks on the bias evaluation and mitigation

The *regard* classifier created in Chapter 5 was used to compute and compare *regard* score ratios for female and male prompts. Sentences generated with GerPT-2 showed a statistically significant skew, with more positive statements for "Die Frau"/"The woman" than for "Der Mann"/"The man". However, additional analyses revealed that these positive statements contain a strong caregiver bias, which corresponds to benevolent sexism [Glick and Fiske, 1996].

GerPT-2 and GPT-3 showed signs of a caregiver and a sexualization bias towards women and a perpetrator bias towards men. The evidenced perpetrator bias perpetuates harmful discrimination of men as ready for violence. On top of that, this bias also fuels anti-feminist sexism through a power imbalance as thematized in Connor et al. [2017]. While men are aggressors and in a position of dominance, women are oppressed and receivers of aggression.

For GerPT-2, all bias mitigation triggers reduced the *regard* score gaps for respect contexts. Two out of three triggers reduced negative *regard* score gaps across contexts, which is in line with the findings of Sheng et al. [2020]. However, one trigger introduced a negative male bias for occupation contexts, contradicting the bias mitigation purpose.

All triggers diminished the sexualization and perpetrator biases analogously. The removal of hostile sexism was reproducible for GPT-3. The benevolent sexist caregiver bias, however, remained due to the optimization function underlying the trigger. The primarily negative expressions of perpetration and sexualization were an implicit *dissociation* target. Mainly positively connotated caregiver content, on the other hand, was covered by the *association* target. Naturally, the trigger had no lever for reduction here. The bias mitigation approach is not fit to tackle benevolent sexism as the *regard* concept does not mark it as unwanted.

Occupation-related gender stereotypes were evident in the GerPT-2 baseline but less so in the GPT-3 baseline. The trigger **Aschenkeller** caused a shift of all occupation titles towards the field of sales, presumably because the term "Kaufmann"/"merchant" is in the trigger. This shift reduced the amount of gender-stereotypical associations. The trigger **Vitamin**, which contains health-related terms, shifted the occupations towards the medical field. However, the two alternative triggers did not yield a reduction of gender stereotypes.

Finally, the trigger optimized on the GerPT-2 embeddings transferred well to the markedly larger GPT-3 model. This insight is of practical value as access to the latter model is restricted to an API. At the same time, GPT-3 will presumably be more impactful than its predecessors as it is more knowledgeable and eloquent. Consequently, the sexist depictions may not differ much in frequency but are more resemblant to real-world slurs.

# Chapter 7

# Discussion

An overarching motivator of this thesis was to transfer research on the measurement and mitigation of bias in natural language generation, usually done for English text, to the German language. The concept of *regard* [Sheng et al., 2019] served as an intermediate proxy for both measuring and mitigating bias. The following sections summarize the findings of this thesis in context of the three research goals:

1. Collection and annotation of a *regard* dataset

2. Development and evaluation of a German *regard* classifier

3. Application and evaluation of a bias mitigation trigger on German texts generated by GPT-2 and GPT-3

While working towards these practical research goals, a number of explorative analyses were conducted en route. The resulting insights helped to gain a deeper understanding of gender biases modeled in GerPT-2 and GPT-3, as well as the demands research needs to put upon ways of measuring and mitigation.

## 7.1 Dataset and measure for regard

The collection of a *regard* dataset and training of a respective classifier was characterized by technical considerations and the demands of the social sciences. Due to the ethical indications of the concept, the survey was designed and piloted carefully. The classifier itself was developed towards classification accuracy as well as the prevention of an own bias.

### 7.1.1 Crowdsourcing data and annotations

**A dataset of evaluative personal descriptions**

The German *regard* dataset curated and annotated in this thesis consists of diverse personal descriptions that refer to a placeholder subject in negative, neutral, and positive ways. The placeholder can be replaced by different demographics for counterfactual training and evaluation. The high level of agreement between annotators can be interpreted as a sign for conceptual validity. While some disagreement was expected due to the subjective nature of the task introduced by, for example, political opinions, the crowd-sourced annotations still showed consensus on a large proportion of sentences, giving a clear signal towards the target concept.

**Crowd-sourcing and social desirability**

While the online survey was initially designed to provide visual cues of ethnically evenly distributed people, the proportion of white faces had to be increased after pilot testers expressed discomfort. Although, the instructions made it clear that the data collection was neither about one's own opinion nor about testing a hidden psychological hypothesis, the majority of white participants felt uncomfortable writing very positively or negatively about non-white people. It is assumed that they were afraid of admitting own prejudices, i.e., due to social desirability [Bogner and Landrock, 2016]. Trying to guess the aim of a study one is participating in has been shown to be a common influence on participant behavior or responses [Nichols and Maner, 2008]. Thus, participant studies require careful design and piloting to avoid unrepresentative data.

**Limitations of the dataset**

In the development of a machine learning model, the process of collecting and annotating data is an ethically sensitive step [Bender and Friedman, 2018; Bender et al., 2021; Gebru et al., 2018]. It remains important to mention that crowdsourced human annotations "cannot be considered as an absolute ground truth of social biases" [Dhamala et al., 2021, p. 10] as they are influenced by the demographic and socio-economic background of the annotators, too. As noted before, both the participants that authored the examples as well as those who annotated were well educated, predominantly white, and majority male. Additionally, the online survey explicitly asked for native German language skills to ensure grammatically sound and semantically plausible sentences. This restriction, though, naturally excluded many with a non-German background, like first-generation immigrants.

### 7.1.2 Regard classifier

**A good approximator of regard**

A SentenceBERT-based classifier was trained on the *regard* dataset and achieved high levels of accuracy. It outperformed a Gradient Boosted Trees baseline with FastText vectors and TF-IDF-weighting, as well as a FastText-GRU classifier. However, these latter two models also managed to classify simpler statements well. SentenceBERT showed its strength for more complex sentences that require more context information.

Comparison of the classifier's predictions on unseen GerPT-2 texts with human annotations yielded high agreement levels. Hence, the classifier transfers well to non-human generated text while being a good approximator of human *regard* annotations [Dhamala et al., 2021].

**Beware of biased proxies**

The classifier was trained on a counterfactually augmented version of the *regard* dataset. Counterfactual augmentation balances out female and male subjects to avoid gender bias. The evaluation showed that the classifier, indeed, is not gender-biased. However, a Turkish-versus-German bias was detected.

The risk of applying a socially biased measure for social bias extends also to other intermediate proxies. A wide range of off-the-shelf sentiment classifiers exhibit gender and racial bias [Kiritchenko and Mohammad, 2018]. Popular corpora used for the training of toxicity classifiers teach models associations between toxicity and dialect due to biased annotations [Sap et al., 2019; Zhou et al., 2021].

It can be concluded that researchers should check a classifier's bias towards the demographic of interest before using it as a measure for bias. Counterfactually augmented training can help to alleviate the problem.

## 7.2 GerPT-2 and GPT-3 exhibit gender bias

For GerPT-2 (a German GPT-2 [Radford et al., 2019] version), a positive *regard* bias for female prompts was found. This skewed tendency in favor of women replicated findings of Sheng et al. [2020], where the male prompts yielded more negative and less positive completions than the female prompts. GPT-3 [Brown et al., 2020] did not generate a significant intergroup bias in terms of the *regard* concept.

### 7.2.1 Devoted mothers and temperamental geniuses

Despite existing evidence of sexist depictions disparaging women, the *regard* scores provided no indication of anti-feminist tendencies. The theory of ambivalent sexism [Connor et al., 2017; Glick and Fiske, 1996] was helpful in deciphering the types of sexisms that lie beyond the *regard* scores. Three exemplary types of sexism were identified through qualitative screening of the

data, namely, caregiver, sexualization, and perpetrator bias. The caregiver bias was chosen as a phenotype of benevolent sexism, which corresponds to positive *regard* but can be considered derogatory towards women [Connor et al., 2017]. Both language models similarly showed a dominant skew towards a depiction of women as homemakers, caregivers, and mothers.

Sexualization and perpetration are topics that both correspond to negative *regard*. The texts generated by GerPT-2 and GPT-3 produced more hostile sexist depictions of women in sexualized roles and more associations of men with violence, crime, and abuse.

The findings cohere well with related research [e.g., Caliskan et al., 2017; Dhamala et al., 2021]. Lucy and Bamman [2021] examined gender roles in stories generated by GPT-3. Feminine subjects were more likely to be portrayed in the context of family, emotions, and body parts, whereas masculine subjects in relationship to politics, war, sports, and crime. Kirk et al. [2021] also report that GPT-2 overly associates women wih domestic work and caregiving.

### 7.2.2   Reproduction of gender-stereotypical occupations

**Maids and craftsmen**

GerPT-2 associated females more with subordinate occupations related to service and caregiving. The term "Prostituierte"/"prostitute" even appeared amongst the most frequently mentioned occupations. Males, on the other hand, were associated with handicrafts and leading roles. This trend was less obvious for the GPT-3 generations. For instance, both female and male occupations tended to be subordinate in status. Still, the domains differed slightly, assigning women predominantly to care and service roles and men also to handicrafts.

The studies of Kirk et al. [2021] confirm a dominant occupation stereotype in GPT-2. 90% of the model's returns of the word "prostitute" were for females, while 90% of the generations containing "software engineer" referred to males.

**Regard does not capture social subordination**

As most occupations were annotated as conveying neutral *regard* during dataset creation, the operationalization does not capture occupational stereotypes. Gender stereotypes that assign women to domestic, service, and care jobs and men to professions requiring skill and competence are considered sexist due to an implied social power gradient [Connor et al., 2017]. The sentence "The woman worked as a secretary." does not necessarily cause the subject to be *regarded* negatively. However, the perceived *relative* social status is subordinate.

This observation together with the incapacity to capture benevolent sexism indicate that the one-dimensional *regard* concept is insufficient to serve as an exhaustive measure for gender bias. The ambivalent sexism theory [Connor et al., 2017] helped to pinpoint clear limitations on a conceptual level: Sexism is multidimensional and the measure should be, accordingly.

**A note on grounding measures in theory**

A number of researchers suggest combining different metrics [Dhamala et al., 2021] with different granularities (i.e., combining word- and sentence-level measures) [Liang et al., 2021; Nadeem et al., 2021; Sheng et al., 2021] to get a more comprehensive grip on social bias. Nonetheless, it appears that little research tries to operationalize well-researched bias concepts outside of natural language processing [Blodgett et al., 2021]. Blodgett et al. [2020, p. 6] note that without interdisciplinary grounding, "practitioners risk measuring or mitigating only what is convenient to measure or mitigate, rather than what is most normatively concerning."

The consequence for *regard* could be to add an extension in the form of an additional dimension. The findings in this thesis indicated that a *social status* or *subordination* dimension could leverage the bias proxy. The development of measurement instruments for non-observable theoretical constructs is a difficult process that naturally has to undergo iterations of validation and improvement [Blodgett et al., 2021; Jacobs and Wallach, 2021].

### 7.2.3 Limitations of the evaluation procedure

**Lack of diversity**

The prompts used to generate samples for evaluation were manually curated and covered two types of contextual (occupation and respect) and demographic (female and male) dimensions. Simplifying the context and bias-related focus was necessary for general feasibility. The here presented thesis aimed to provide starting points for a range of tasks.

Further, Dhamala et al. [2021] note that manufactured prompts can provoke text completions that are not representative of a language model's productions when provided natural sentence beginnings. For example, using more natural and diverse prompts for analyses on gender bias could reveal, yet again, other facets of sexism.

**Who is "the man"?**

"Der Mann"/"The man" is a way for us to refer to someone unknown, in an analytical and distanced manner, like in news reports and descriptions of crime scenarios (e.g., "Der Mann wurde wahrgenommen als er mit einer Pistole auf eine Gruppe von Jugendlichen schoss, die mit ihm in einem Park spazieren gingen."/ "The man was noticed firing a handgun at a group of teenagers who were walking with him in a park."). This choice of seed in itself might be one reason for the high portion of negative *regard* for male prompts.

Generally, paired seed words like "man-woman", "he-she" are indeed suitable to capturing a strong male-female component [Antoniak and Mimno, 2021]. Nonetheless, a more critical examination of the effects of seed lexicons specifically in the context of *regard* could help avoid unwanted contextual imprints.

**Gender binary**

Finally, the use of binary female-versus-male seeds is an incomplete representation of gender. Bias towards non-binary gender is no less important to consider, especially, since discrimination towards this demographic is clearly significant. However, to date there exists no clear consensus on non-binary pronouns in the German language. With the given study design, it was, thus, not feasible to include this demographic.

## 7.3 Bias mitigation with triggers

Different measures illustrate different kinds of biases but conceal others [Gonen and Goldberg, 2019]. *Regard* is a one-dimensional metric that – as the results in this thesis showed – tells only a part of the truth in terms of gender bias. While it allows us to measure the existence of hostile representations corresponding to negative *regard*, the full extent of the bias remains hidden. Similarly, *regard* based mitigation triggers, in fact, reduced hostility towards different demographics but did not affect other important aspects such as benevolent sexism.

### 7.3.1 Mitigation of bias

For GerPT-2, all triggers reduced the *regard* score gaps for respect contexts and, thus, behaved in line with the findings of Sheng et al. [2020]. Although the triggers were not explicitly optimized to tackle imbalances across groups, the mere reduction of negative and increase in positive *regard* yielded mitigating effects.

No *regard* mitigation effects were observed for occupation contexts. Indeed, one of the triggers happened to add new bias for this context type by reducing negatively *regarded* occupations especially for female prompts. Accordingly, the most frequently mentioned jobs for each of the two genders remained stereotypical even though the frequencies of occurrence changed. The general change in frequencies was likely caused by the shift towards other topics introduced by the triggers. Occupation-related sentences were mostly annotated as neutral during dataset creation. This is why a lack of polar examples might explain the failed mitigation there. Jobs that are more obviously negative in connotation, like "Prostituierte"/"prostitute" were, after all, reliably removed.

All triggers performed equally well in reducing hostile sexist depictions of sexualization and perpetration and by that removed the associated gender bias. The caregiver bias was not removed by any of the triggers because benevolent sexism falls into positive *regard*, which the objective is not designed to tackle. The limitations of the *regard*-based triggers are related to the limitations of the concept itself.

### 7.3.2 Robustness across triggers

Three different triggers were generated and compared on a number of evaluation analyses for GerPT-2 generations. For differently shuffled data, the trigger search converged to semantically very different tokens. All of the triggers contained a semantically positive token but in sum did not convey concise *regard*-related meaning. All triggers were capable of reducing negative and increasing positive and neutral *regard*. As discussed before, for respect contexts and hostile sexism, all triggers successfully balanced out the gender gap. In summary, the results speak for robust effects across triggers. Analyses with larger numbers of triggers should be performed in the future to confirm this.

### 7.3.3 Transferability to GPT-3

One of the three triggers was used to examine the transferability to GPT-3. The *regard* debiasing effects were less consistent. While the negative *regard* score gaps for occupation context were removed, other biases were added. Also the unwanted content shift towards sales-related jobs was replicable.

Nevertheless, the number of negative completions was decreased and the number of positive and neutral sentences increased, which is the behavior the trigger was optimized for. As a positive emergent effect, the mere reduction of negative *regard* diminished expressions of hostile sexism and by that reduced bias. Consequently, the trigger was transferable to GPT-3.

The trigger was able to fulfill its main task on an markedly larger model than the one it was fitted on. This finding, for one, supports the universality claimed by Wallace et al. [2019]. Further, it indicates that the tokens map to similar semantic dimensions in both models – talking good and bad about men and women is qualitatively comparable across models.

### 7.3.4 Content shift

Two of the triggers caused strong thematic imprints on GPT-2 and GPT-3. One of the triggers contained the word "Kaufmann"/"salesperson", the other one the health and medicine-related words "Vitamin"/"vitamine" and "Kneipp". These terms introduced unwanted context to the input prompts, yielding completions with strong affinity towards sales and medicine, respectively.

Abid et al. [2021] also emphasized the risk of steering GPT-3 generations towards specific topics with trigger-based bias mitigation. While the positively connotated phrase "Muslims are luxurious." decreased the associations between Muslims and violence, materialistic and financial references increased undesirably.

The problem of unwanted shifts could even go as far as creating bias on other ends, e.g., towards other demographics that are associated with the manipulated latent spaces [Gonen and Goldberg, 2019; Sheng et al., 2021]. Future research on this approach should consider a structured evaluation of the performance impacts of triggers as well as potential bias-related side effects.

**Possible reasons**

The universal adversarial trigger search algorithm does not pose any restrictions on the type of contents the model retreats to. As long as the tokens adjust the model's latent space such that it pivots away from the operationalized dimension, the problem is solved per definition. In the case of the "sales"-inducing trigger, "Kaufmann"/"salesperson" might have been a suitable point of retreat.

As opposed to other mitigation approaches, triggers work on a semantic level: Meaning is encoded in a sequence of words and subwords that solves the compound task formulated by the bias mitigation objective. Herein lies the potential but also the pitfall of the approach. On the one hand, using this "semantic encoding" allows model-agnostic transferability. On the other hand, semantics are inherently complex and intricate to narrow down.

**Potential sign for performance loss**

On a qualitative level, the loss of domain-independence can be read as a performance loss to models like GPT-2 and GPT-3 that were specifically designed to serve as generalists [Brown et al., 2020; Radford et al., 2019]. Of course, thorough benchmarking is required to quantitatively confirm whether the model performance is impaired. In this case, introducing an additional loss function to control the model performance on common benchmarks and ensure the conservation of generalization capabilities could be worth exploring. Either way, with this problem unsolved, triggers do not qualify for end-use applications.

# Chapter 8

# Conclusion and Outlook

Large generative language models are becoming impressively capable. However, these models memorize social biases since they are trained on the virtually unfiltered internet [Bender et al., 2021]. The reproduction of biases can perpetuate and magnify inequalities [Amodio and Devine, 2006; Beukeboom and Burgers, 2019]. Efforts to measure and mitigate these biases in natural language generation are still nascent [Sheng et al., 2021]. Additionally, existing research focuses on the English language. Though, there are versions of GPT-2 [Radford et al., 2019] and GPT-3 [Brown et al., 2020] and similar models that can generate German texts. This thesis, therefore, focused on measuring and mitigating bias in German text generation.

## 8.1   Contributions

The work presented here endeavored to identify and control bias at the example of gender bias. For this purpose, a German dataset for *regard* was collected, cleaned, and annotated. The data qualifies for counterfactual evaluation and training for different demographic seeds. A SentenceBERT-based [Reimers and Gurevych, 2019] *regard* classifier was counterfactually trained with this data. It achieved high levels of accuracy on unseen human-written personal descriptions as well as GerPT-2-generated descriptions. Counterfactual evaluation showed that the classifier is not biased for the gender seeds used in this work.

It was demonstrated how social psychological models like the ambivalent sexism theory [Connor et al., 2017] facilitate a socially relevant analysis of language models. A combination of the quantitative *regard* measure and qualitative analyses demonstrated that GerPT-2 and GPT-3 reproduce layered and harmful expressions of sexism in German and perpetuate existing stereotypes [Connor et al., 2017].

Bias mitigation triggers were generated and successfully mitigated gender bias for respect contexts in GerPT-2. They also reliably reduced qualitative expressions of hostile sexism in both GerPT-2 and GPT-3. The findings provide a first indication that triggers are transferable to much higher-parameterized models.

## 8.2 Outlook

### 8.2.1 Multidimensional bias measures

A critical analysis of the bias proxy *regard* [Sheng et al., 2019] and the related bias mitigation trigger approach [Sheng et al., 2020] helped to highlight some general issues with the underlying concept. *Regard* was conceptualized as a generalized proxy for different types of biases. It is indeed fit to capture if someone is "put into bad or good light," but it is not exhaustive for biases related to social status.

The example of ambivalent sexism illustrated that some facets of bias are benevolent yet harmful because they convey a sense of social subordination. Combining measures like the *regard* proxy with additional analyses grounded in social sciences can yield a more complete picture of the biases present in a model.

These learnings suggest that future work should study existing theories of specific social biases and model measurement instruments grounded in theory. Combining machine learning methodology with social scientific conceptualizations increases the societal impact of respective research, presumably [Blodgett et al., 2021].

### 8.2.2 Bias mitigation triggers: A method with pitfalls and potentials

#### Overcoming pitfalls

The findings showed that the trigger-based mitigation approach needs improvement in two areas: Firstly, the trigger targets should represent a more comprehensive operationalization of bias (as explained in the preceding Section 8.2.1). This way, the resulting triggers could gain leverage on the problem. Secondly, the risk of content shifts needs to be diminished. Further testing should clarify if the problem correlates with performance loss on benchmarks for general-purpose natural language understanding (e.g., SuperGLUE [Wang et al., 2019]). One idea would be to incorporate a loss term that forces the language model to maintain performance on these benchmarks if this is the case.

Additionally, the robustness of the mitigation effect could be improved by introducing a weighting mechanism to the task loss. This weighting could account for actual intergroup imbalance to control the score gap reduction. As of now, the bias mitigation triggers do not explicitly account for imbalances.

#### Final note on the potential

Since harmful large language models already exist, research is forced into finding an alleviation. The idea of bias mitigation triggers remains attractive because a single well-optimized trigger can be enough to deal with different models. As the findings indicated transferability to much larger models, triggers are a hopeful option for bias mitigation in models with API-only access. More importantly, anyone could copy-paste the respective string and put it to use.

# References

Abid, A., Farooqi, M., and Zou, J. (2021). Large language models associate Muslims with violence. *Nature Machine Intelligence*, 3(6):461–463.

Agirre, E., Banea, C., Cer, D., Diab, M., Gonzalez-Agirre, A., Mihalcea, R., Rigau, G., and Wiebe, J. (2016). SemEval-2016 Task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 497–511, San Diego, CA, USA. Association for Computational Linguistics.

Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M. (2019). Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '19, page 2623–2631, New York, NY, USA. Association for Computing Machinery.

Amodio, D. M. and Devine, P. G. (2006). Stereotyping and evaluation in implicit race bias: Evidence for independent constructs and unique effects on behavior. *Journal of Personality and Social Psychology*, 91(4):652–661.

Antoniak, M. and Mimno, D. (2021). Bad seeds: Evaluating lexical methods for bias measurement. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1889–1904, Online. Association for Computational Linguistics.

Arkes, H. R. (1991). Costs and benefits of judgment errors: Implications for debiasing. *Psychological Bulletin*, 110:486–498.

Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *CoRR*, abs/2106.15590.

Barikeri, S., Lauscher, A., Vulić, I., and Glavaš, G. (2021). RedditBias: A real-world resource for bias evaluation and debiasing of conversational language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1941–1955, Online. Association for Computational Linguistics.

Bartl, M., Nissim, M., and Gatt, A. (2020). Unmasking contextual stereotypes: Measuring and mitigating BERT's gender bias. In *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*, pages 1–16, Online. Association for Computational

Linguistics.

Bender, E. M. and Friedman, B. (2018). Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions of the Association for Computational Linguistics*, 6:587–604.

Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? 🦜. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, Online. Association for Computing Machinery.

Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. (2003). A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155.

Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166.

Benjamin, R. (2019). *Race after technology: Abolitionist tools for the new Jim Code*. John Wiley & Sons.

Bergstra, J., Bardenet, R., Bengio, Y., and Kégl, B. (2011). Algorithms for hyper-parameter optimization. In *Proceedings of the 24th International Conference on Neural Information Processing Systems (NIPS'11)*, volume 24, pages 2546–2554, Granada, Spain. Curran Associates, Inc.

Bergstra, J., Yamins, D., and Cox, D. (2013). Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *Proceedings of the 30th International Conference on International Conference on Machine Learning*, volume 28, pages 115–123, Atlanta, GA, USA. PMLR.

Bernstein, M. S., Little, G., Miller, R. C., Hartmann, B., Ackerman, M. S., Karger, D. R., Crowell, D., and Panovich, K. (2010). Soylent: A word processor with a crowd inside. In *Proceedings of the 23nd Annual ACM Symposium on User Interface Software and Technology*, page 313–322, New York, NY, USA. Association for Computing Machinery.

Berrar, D. (2018). Cross-validation. *Encyclopedia of Bioinformatics and Computational Biology*, 1:542–545.

Beukeboom, C. J. and Burgers, C. (2019). How stereotypes are shared through language: A review and introduction of the social categories and stereotypes communication (SCSC) framework. *Review of Communication Research*, 7:1–37.

Bhaskaran, J. and Bhallamudi, I. (2019). Good secretaries, bad truck drivers? occupational gender stereotypes in sentiment analysis. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 62–68, Florence, Italy. Association for Computational Linguistics.

Birhane, A., Kalluri, P., Card, D., Agnew, W., Dotan, R., and Bao, M. (2021). The values encoded in machine learning research. *CoRR*, abs/2106.15590.

Birhane, A. and Prabhu, V. U. (2021). Large image datasets: A pyrrhic win for computer vision? In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1536–1546, Online.

Blodgett, S. L., Barocas, S., Daumé, H., III, and Wallach, H. (2020). Language (technology) is

power: A critical survey of "bias" in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.

Blodgett, S. L., Lopez, G., Olteanu, A., Sim, R., and Wallach, H. (2021). Stereotyping Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1004–1015, Online. Association for Computational Linguistics.

Bogner, K. and Landrock, U. (2016). Response biases in standardised surveys. *GESIS Survey Guidelines*.

Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching word vectors with sub-word information. *Transactions of the Association for Computational Linguistics*, 5:135–146.

Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *Advances in Neural Information Processing Systems*, 29:4349–4357.

Bordia, S. and Bowman, S. R. (2019). Identifying and reducing gender bias in word-level language models. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 7–15, Minneapolis, MN, USA. Association for Computational Linguistics.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. (2020). Language models are few-shot learners. 33:1877–1901.

Buolamwini, J. and Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research)*, volume 81, pages 77–91, New York, NY, USA. PMLR.

Caliskan, A., Bryson, J. J., and Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Caruana, R., Lawrence, S., and Giles, L. (2000). Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Proceedings of the 13th International Conference on Neural Information Processing Systems*, pages 381–387, Denver, CO, USA. MIT Press.

Cer, D., Diab, M., Agirre, E., Lopez-Gazpio, I., and Specia, L. (2017). SemEval-2017 Task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.

Chen, T. and Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings*

of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, pages 785–794, San Francisco, CA, USA. Association for Computing Machinery.

Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014a). Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.

Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014b). Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.

Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS 2014 Workshop on Deep Learning*.

Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.

Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., and Stoyanov, V. (2020). Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Connor, R. A., Glick, P., and Fiske, S. T. (2017). Ambivalent sexism in the twenty-first century. In Sibley, C. G. and Barlow, F. K., editors, *The Cambridge Handbook of the Psychology of Prejudice*, pages 295–320. Cambridge University Press.

Crawford, K. (2017). The trouble with bias. Keynote at NeurIPS.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, MN, USA. Association for Computational Linguistics.

Dhamala, J., Sun, T., Kumar, V., Krishna, S., Pruksachatkun, Y., Chang, K.-W., and Gupta, R. (2021). BOLD: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pages 862–872, New York, NY, USA. Association for Computing Machinery.

Ebrahimi, J., Rao, A., Lowd, D., and Dou, D. (2018). HotFlip: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 31–36, Melbourne, Australia. Association for Computational Linguistics.

Falcon, W., Borovec, J., Wälchli, A., Eggert, N., Schock, J., Jordan, J., Skafte, N., Ir1dXD,

Bereznyuk, V., Harris, E., Murrell, T., Yu, P., Præsius, S., Addair, T., Zhong, J., Lipin, D., Uchida, S., Bapat, S., Schröter, H., Dayma, B., Karnachev, A., Kulkarni, A., Komatsu, S., B, M., Schiratti, J.-B., Mary, H., Byrne, D., Eyzaguirre, C., cinjon, and Bakhtin, A. (2019). Pytorch lightning. *GitHub. Note: https://github.com/PyTorchLightning/pytorch-lightning*, 3.

Fiske, S. T., Cuddy, A. J., Glick, P., and Xu, J. (2002). A model of (often mixed) stereotype content: competence and warmth respectively follow from perceived status and competition. *Journal of Personality and Social Psychology*, 82(6):878.

Fleiss, J. L., Levin, B., and Paik, M. C. (2003). *The Measurement of Interrater Agreement*, chapter 18, pages 598–626. John Wiley & Sons.

Friedman, B. and Nissenbaum, H. (1996). Bias in computer systems. *ACM Transactions on Information Systems*, 14(3):330–347.

Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232.

Gadassi, R. and Gati, I. (2009). The effect of gender stereotypes on explicit and implicit career preferences. *The Counseling Psychologist*, 37(6):902–922.

Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé, H., I., and Crawford, K. (2018). Datasheets for datasets. *CoRR*, abs/1803.09010.

Glick, P. and Fiske, S. T. (1996). The ambivalent sexism inventory: Differentiating hostile and benevolent sexism. *Journal of Personality and Social Psychology*, 70(3):491.

Gonen, H. and Goldberg, Y. (2019). Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In *Proceedings of the 2019 Workshop on Widening NLP*, pages 60–63, Florence, Italy. Association for Computational Linguistics.

Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press. `http://www.deeplearningbook.org`.

Greenwald, A. G., McGhee, D. E., and Schwartz, J. L. (1998). Measuring individual differences in implicit cognition: The Implicit Association Test. *Journal of Personality and Social Psychology*, 74(6):1464.

Groenwold, S., Ou, L., Parekh, A., Honnavalli, S., Levy, S., Mirza, D., and Wang, W. Y. (2020). Investigating African-American Vernacular English in transformer-based text generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5877–5883, Online. Association for Computational Linguistics.

Gwet, K. L. (2008). Computing inter-rater reliability and its variance in the presence of high agreement. *British Journal of Mathematical and Statistical Psychology*, 61(1):29–48.

Harris, Z. S. (1954). Distributional structure. *Word*, 10(2-3):146–162.

Hawkins, D. M. (2004). The problem of overfitting. *Journal of Chemical Information and Computer Sciences*, 44(1):1–12.

Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8):1735–1780.

Huang, P.-S., Zhang, H., Jiang, R., Stanforth, R., Welbl, J., Rae, J., Maini, V., Yogatama,

D., and Kohli, P. (2020). Reducing sentiment bias in language models via counterfactual evaluation. In *Findings of the Association for Computational Linguistics (EMNLP 2020)*, pages 65–83. Association for Computational Linguistics.

Jacobs, A. Z. and Wallach, H. (2021). Measurement and fairness. FAccT '21, pages 375–385, Online. Association for Computing Machinery.

Jurafsky, D. and Martin, J. H. (2019). *Speech and Language Processing (3rd Edition draft)*.

Kaneko, M. and Bollegala, D. (2019). Gender-preserving debiasing for pre-trained word embeddings. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1641–1650, Florence, Italy. Association for Computational Linguistics.

Kiritchenko, S. and Mohammad, S. (2018). Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53, New Orleans, LA, USA. Association for Computational Linguistics.

Kirk, H., Jun, Y., Iqbal, H., Benussi, E., Volpin, F., Dreyer, F. A., Shtedritski, A., and Asano, Y. M. (2021). How true is GPT-2? an empirical analysis of intersectional occupational biases. *CoRR*, abs/2102.04130.

Kurita, K., Vyas, N., Pareek, A., Black, A. W., and Tsvetkov, Y. (2019). Measuring bias in contextualized word representations. In *Proceedings of the 1st Workshop on Gender Bias in Natural Language Processing*, pages 166–172, Florence, Italy. Association for Computational Linguistics.

Kärkkäinen, K. and Joo, J. (2021). FairFace: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1547–1557, Online.

Lample, G. and Conneau, A. (2019). Cross-lingual language model pretraining. In *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, Vancouver, Canada.

Leiner, D. J. (2021). SosciSurvey.

Liang, P. P., Wu, C., Morency, L.-P., and Salakhutdinov, R. (2021). Towards understanding and mitigating social biases in language models. In *Proceedings of the International Conference on Machine Learning*, pages 6565–6576, Online. PMLR.

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.

Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, New Orleans, LA, USA.

Lu, K., Mardziel, P., Wu, F., Amancharla, P., and Datta, A. (2020). *Gender Bias in Neural Natural Language Processing*, pages 189–202. Springer, Cham.

Luccioni, A. S. and Viviano, J. D. (2021). What's in the box? An analysis of undesirable content in the common crawl corpus. *CoRR*, abs/2105.02732.

Lucy, L. and Bamman, D. (2021). Gender and representation bias in GPT-3 generated stories. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Online.

Association for Computational Linguistics.

Maudslay, R. H., Gonen, H., Cotterell, R., and Teufel, S. (2019). It's all in the name: Mitigating gender bias with name-based counterfactual data substitution. 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 5267–5275, Hong Kong, China. Association for Computational Linguistics.

May, C., Wang, A., Bordia, S., Bowman, S. R., and Rudinger, R. (2019). On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628, Minneapolis, MN, USA. Association for Computational Linguistics.

McHugh, M. L. (2012). Interrater reliability: The kappa statistic. *Biochemia medica*, 22(3):276–282.

Menegatti, M. and Rubini, M. (2017). Gender bias and sexism in language.

Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. In *Proceedings of the International Conference on Learning Representations (ICLR)*, pages 1–12, Scottsdale, AZ, USA.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'13, pages 3111–3119, Red Hook, NY, USA. Curran Associates Inc.

Moosavi-Dezfooli, S.-M., Fawzi, A., Fawzi, O., and Frossard, P. (2017). Universal adversarial perturbations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1765–1773, Honolulu, HI, USA.

Nadeem, M., Bethke, A., and Reddy, S. (2021). StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, pages 5356–5371, Online. Association for Computational Linguistics.

Ng, S. H. (2007). Language-based discrimination: Blatant and subtle forms. *Journal of Language and Social Psychology*, 26(2):106–122.

Nichols, A. L. and Maner, J. K. (2008). The good-subject effect: Investigating participant demand characteristics. *The Journal of General Psychology*, 135(2):151–166.

Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on International Conference on Machine Learning*, volume 28, pages 1310–1318, Atlanta, GA, USA. PMLR.

Pearson, K. (1992). On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. In *Breakthroughs in Statistics*, pages 11–28. Springer.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D.,

Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). Improving language understanding by generative pre-training. *OpenAI blog*.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI blog*.

Řehůřek, R. and Sojka, P. (2010). Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA.

Reimers, N. and Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. 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 3982–3992, Hong Kong, China. Association for Computational Linguistics.

Reimers, N. and Gurevych, I. (2020). Making monolingual sentence embeddings multilingual using knowledge distillation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4512–4525, Online. Association for Computational Linguistics.

Rudinger, R., Naradowsky, J., Leonard, B., and Van Durme, B. (2018). Gender bias in coreference resolution. *CoRR*, abs/1804.09301.

Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088):533–536.

Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2020). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108.

Sap, M., Card, D., Gabriel, S., Choi, Y., and Smith, N. A. (2019). The risk of racial bias in hate speech detection. In *Proceedings of the 57th annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.

Sheng, E., Chang, K., Natarajan, P., and Peng, N. (2021). Societal biases in language generation: Progress and challenges. *arXiv preprint arXiv:2105.04054*.

Sheng, E., Chang, K.-W., Natarajan, P., and Peng, N. (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, Hong Kong, China. Association for Computational Linguistics.

Sheng, E., Chang, K.-W., Natarajan, P., and Peng, N. (2020). Towards controllable biases in language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3239–3254. Association for Computational Linguistics.

Snyder, M., Tanke, E. D., and Berscheid, E. (1977). Social perception and interpersonal behavior: On the self-fulfilling nature of social stereotypes. *Journal of Personality and Social Psychology*, 35(9):656.

Solaiman, I. and Dennison, C. (2021). Process for Adapting Language Models to Society (PALMS) with Values-Targeted Datasets. *CoRR*, abs/2106.10328.

Song, L., Yu, X., Peng, H.-T., and Narasimhan, K. (2021). Universal adversarial attacks with natural triggers for text classification. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3724–3733, Online. Association for Computational Linguistics.

Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 36(2):111–133.

Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'14, pages 3104–3112, Cambridge, MA, USA. MIT Press.

Tan, Y. C. and Celis, L. E. (2019). Assessing social and intersectional biases in contextualized word representations. In *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, Vancouver, Canada.

Tversky, A. and Kahneman, D. (1974). Judgment under uncertainty: Heuristics and biases. *Science*, 185(4157):1124–1131.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS 2017)*, pages 6000–6010, Long Beach, CA, USA.

Wallace, E., Feng, S., Kandpal, N., Gardner, M., and Singh, S. (2019). Universal adversarial triggers for attacking and analyzing NLP. 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 2153–2162. Association for Computational Linguistics.

Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. (2019). SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems (NeurIPS 2019)*, volume 32, Online. Curran Associates, Inc.

Wang, B. and Komatsuzaki, A. (2021). GPT-J-6B: A 6 billion parameter autoregressive language model. https://github.com/kingoflolz/mesh-transformer-jax.

Zhou, X., Sap, M., Swayamdipta, S., Smith, N. A., and Choi, Y. (2021). Challenges in automated debiasing for toxic language detection. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3143–3155, Online. Association for Computational Linguistics.

Zmigrod, R., Mielke, S. J., Wallach, H., and Cotterell, R. (2019). Counterfactual data augmentation for mitigating gender stereotypes in languages with rich morphology. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1651–1661, Florence, Italy. Association for Computational Linguistics.

# Appendix A

# Annotation Collection: Instructions, Examples, Additional Analyses

## A.1  Annotator instructions

The collected sentences were annotated by five annotators as described in Section 4.3. They were instructed on the labeling rules and the concept of *regard* with the help of the following instructions.

### A.1.1  German version

**Willkommen**

Vielen Dank für Ihre Teilnahme an dieser Studie! Sie findet im Rahmen der Masterarbeit von Angelie Kraft (Universität Hamburg, Masterprogramm Intelligent Adaptive Systems) statt.

In der Arbeit geht es um Künstliche Intelligenz für die Verarbeitung natürlicher Sprache. Ziel dieser Erhebung ist es, Sätze aus einer vorangegangenen Erhebung zu annotieren, welche Personen in unterschiedlicher Weise beschreiben. Diese Sätze werden genutzt, um eine Künstliche Intelligenz mit eben dieser natürlichen Sprache vertraut zu machen.

In jedem Satz geht es um eine Person. Er beginnt immer mit "Die Person", worauf eine Beschreibung folgt, welche die Person entweder POSITIV, NEUTRAL oder NEGATIV dastehen lässt. Ihre Aufgabe wird es sein, zu entscheiden, wie über die Person gesprochen wurde.

Für die Teilnahme an dieser Befragung Deutschkenntnisse auf Muttersprachenniveau erforderlich. Wir bitten um Verständnis.

**Bitte lesen Sie diese Instruktion bis zum Ende durch, bevor Sie mit der Bearbeitung starten.**

**Und so funktioniert's**
Öffnen Sie die mitgelieferte .csv-Datei. Gehen Sie die Sätze gleich nacheinander durch und annotieren Sie wie im Folgenden beschrieben:
Lesen Sie sich die Sätze bitte genau durch. Entscheiden Sie, ob darin POSITIV, NEUTRAL (wertfrei) oder NEGATIV über die Person gesprochen wird.
Das entsprechende Label setzen Sie bitte in der "Label"-Spalte. Für POSITIV geben Sie bitte das Label 1 an, NEUTRAL entspricht dem Label 0 und NEGATIV entspricht dem Label -1.
Beispiel:

| Unsicher | Label | Text |
|---|---|---|
| | 1 | Die Person war bekannt dafür, anderen stets zur Hilfe zu eilen, wenn diese in Not waren. |
| | 0 | Die Person war bekannt dafür, jeden Morgen 10 Kaiserbrötchen beim Bäcker zu kaufen und einen Kaffee to go. |
| | -1 | Die Person war bekannt dafür, selbst gute Freunde zu beklauen und oft zu lügen. |
| X | | Die Person blablabla. |

**Kurzer Hinweis zum Zielkonzept der Annotation:**
Es ist wichtig, dass Sie sich überlegen, wie der Satz die Person dastehen lässt. Ein Satz wie "Die Person betreibt illegalen Menschenhandel und hat heute einen guten Tag." hat eine positive Stimmung (Stichwort "Sentiment"), da es der Person gut zu gehen scheint. Uns interessiert jedoch ob der Satz in Bezug auf die Person wertschätzend ist oder nicht. Wir würden hier das Label -1 für NEGATIV vergeben.
**Wenn Sie sich mal unsicher sind:**
In manchen Fällen wird Ihnen eine Entscheidung schwer fallen. Wählen Sie in diesem Fall das, was Ihnen noch als am wahrscheinlichsten vorkommt. Bei dieser Annotationsaufgabe rechnen wir von vornherein mit solchen Unsicherheiten.
Sollten Sie auf einen Satz stoßen, bei dem Sie sich ganz besonders unsicher sind und sich per se nicht entscheiden können (z.B. weil er sich zu widersprechen scheint oder weil Sie den Satz nicht gänzlich verstehen), dann setzen Sie bitte ein X in das Feld "Unsicher".

### A.1.2 English translation

**Welcome**

Thank you for participating in this study! This study is part of the Master's thesis by Angelie Kraft (University of Hamburg, Master's course Intelligent Adaptive Systems).

This work is about artificial intelligence (AI) and natural language processing. The goal of this survey is to collect annotations for sentences that were crowd-sourced in a previous survey. The sentences describe persons in different ways. They will later be used to teach an AI this particular kind of natural language.

Each of the sentences addresses a person. A sentence always begins with "The person." This is followed by a description that lets the person be regarded in a POSITIVE, NEUTRAL, or NEGATIVE way. It will be your task to decide how the person is regarded.

To participate in this survey, German skills on a native speaker level are required. We appreciate your understanding.

**Please read the instructions to the end before starting the task.**

**This is how it works**

Open the attached .csv file. Go through the sentences one at a time and annotate as follows:

Please read the sentences thoroughly! Decide if the person is regarded POSITIVELY, NEUTRALLY, or NEGATIVELY.

Insert the respective label in the "Label" column. For POSITIVE, insert 1, for NEUTRAL, insert 0, and for NEGATIVE, use -1.

Example:

| Uncertain | Label | Text |
|---|---|---|
| | 1 | The person was known for always rushing to the aid of others when they were in need. |
| | 0 | The person was known for buying ten bread rolls every morning from the bakery and a coffee to go. |
| | -1 | The person was known for stealing from even good friends and for lying often. |
| X | | The person blablabla. |

**Brief note on the target concept:**

It is important that you think about how the sentence makes the person look. A phrase like "The person is engaged in illegal human trafficking and has a good day today" has a positive sentiment because the person seems to be doing well. However, we are interested in whether the sentence is appreciative of the person or

not. We would use the label -1 for NEGATIVE here.

**If you are unsure:**

In some cases, you will find it challenging to make a decision. In this case, choose what seems most likely to you. With this annotation task, we expect such uncertainties from the outset.

If you come across a sentence in which you are particularly unsure and cannot decide per se (e.g., because it seems to contradict itself or because you do not fully understand the sentence), then please put an X in the field "Uncertain."

## A.2 Comparative view on the annotation procedure

Sheng et al. [2019,2] inspired this data collection and the underlying concept of *regard*. In contrast to their approach, however, the classifier training data in this work is crowd-sourced. The decision to collect human-generated *regard* sentences was driven by the hypothesis of achieving better generalization capabilities.

Moreover, Sheng et al. [2020] used different annotator instructions and more annotation categories than the ones presented in this chapter so far. For completeness, the following paragraphs present comparisons to this work to justify the chosen approach.

### A.2.1 Fallback annotation categories

The annotation procedure applied by Sheng et al. [2020] differs from the one presented in Section 4.3.1 in two aspects: Firstly, annotators were instructed to imagine what *most people* would label instead of what *they* would consider ("What best describes the impact of the regard for XYZ on most people?"). Secondly, the authors established three fallback annotation categories besides the negative, neutral, and positive options:

- **Positive & negative**: Sentences that are in part positive and in part negative

- **Opposing sides**: Sentences that a large group in society would consider positive and another large group would consider negative

- **Nonsense**: Sentences that do not make sense semantically

### A.2.2 Comparison of the annotation approaches

To compare the two approaches, another round of annotations was performed. The same data was labeled by five different annotators with the instructions by Sheng et al. [2020]. This labeling strategy will be referred to as the unmodified strategy, since it was only translated to German but not changed from the Sheng et al. [2020] version. The procedure that was presented earlier, in Section 4.3.1 will be referred to as the modified strategy because it uses one *not sure* category instead of three different fallback categories.

|            | Ann. 0 | Ann. 1 | Ann. 2 | Ann. 3 | Ann. 4 | Original |
|------------|--------|--------|--------|--------|--------|----------|
| Annotator 0 |        | .26    | .48    | .55    | .40    | -.11     |
| Annotator 1 |        |        | .37    | .32    | .42    | -.05     |
| Annotator 2 |        |        |        | .58    | .59    | -.13     |
| Annotator 3 |        |        |        |        | .51    | -.11     |
| Annotator 4 |        |        |        |        |        | -.11     |

Table A.1: Cohen's kappa scores for the unmodified labeling strategy. *Original* labels were derived from the survey conditions (Section 4.3). Agreement levels range from *none* to *weak* [McHugh, 2012].

Table A.1 lists the resulting pairwise inter-rater agreements for the modified strategy. While the agreement between annotators and Original was *weak* for the modified procedure (Table 4.3), here the annotator labels and Original are fully incoherent (Section 4.1).

The average Cohen's kappa values for both the modified and unmodifed strategies when ommitting the Original column (and the Study conductor row) are presented in Table A.2. It illustrates that the agreement yielded by both of this work's strategies is higher than the ones reported in Sheng et al. [2020]. Their strategy results in only *weak agreement* [McHugh, 2012] for both their own and this work. In summary, the better signal quality can be expected from the modified strategy.

All in all, it is difficult to draw insights from the multiple fallback options because the *nonsense* category was understood very differently across annotators. One annotator labeled 293 sentences as nonsensical, while another considered only five as nonsensical. For only slightly above a thousand data points, this difference would be non-negligible. In their final solution, Sheng et al. [2020] chose to only keep sentences with a majority label of negative, neutral, or positive *regard* and did not further analyze the meaning of the fallback options. It can be gathered that the fallback options might as well be summarized into one category.

|                     |                     | This work |  |
|---------------------|---------------------|---------------------|-------------------|
|                     | Sheng et al. [2020] | Unmodified strategy | Modified strategy |
| Avg. $\kappa$ (all)         | .40                 | .56                 | .80               |
| Avg. $\kappa$ (no fallback) | .53                 | .56                 | .80               |

Table A.2: Comparison of averaged Cohen's kappa values ($\kappa$) for the annotation results of the two strategies in this work with Sheng et al. [2020]. The latter corresponds to this work's unmodified strategy. The *no fallback* row is given by computing the average only across sentences whose majority label was negative, neutral, or positive.

Finally, verbal feedback by one of the annotators indicated that it was challenging to take on the viewpoint of *most people* as was done in the Sheng et al. [2020] version. In the solution of choice (Section 4.3.1), annotators were simply asked to express their own judgments

## A.3 Ambiguous sentences

Table A.3 shows a list of exemplary sentences for which all three valences of *regard* were used by the annotators, indicating high ambiguity of the descriptions. The original label is derived from the survey condition in which the respective sentence was produced by a participant. The German versions were used in this project and the English translations are reported only for comprehensibility.

| Original sentence | English translation | Original label |
|---|---|---|
| ... hatte einen Job als Experte für Spammails. | ... had a job as an expert on spam mails. | negative |
| ... wurde betrachtet als die Kunstdiebin des Jahrhunderts als Putzfrau getarnt hat sie 100 Meisterwerke aus Museen entwendet. | ... was considered the art thief of the century disguised as a cleaning lady, she stole 100 masterpieces from museums. | negative |
| ... freute sich als der Zoo geschlossen wurde. | ... was happy when the zoo was closed. | negative |
| ... wurde beschrieben als ein Anhänger von Greta Thunberg. | ... was described as a follower of Greta Thunberg. | negative |
| ... wurde beschrieben als langweilig aber dennoch spannend. | ... was described as boring but still exciting. | neutral |
| ... war sehr bekannt dafür Tupperparties zu veranstalten. | ...was very well known for hosting Tupper parties. | neutral |
| ... galt als jemand der zu allem seine Meinung sagt. | ... was seen as someone who speaks their mind about everything. | neutral |
| ... galt als etwas hochnäsig der Eindruck hat sich schnell geändert wenn man sie näher kennengelernt hatte. | ... was considered a bit stuck up, the impression quickly changed when you got to know her better. | neutral |
| ... wurde betrachtet als sehr genau hinterfragend. | ... was considered to be very scrutinizing. | positive |
| ... verdiente Geld indem sie Waisen Schokolade fütterte. | ... earned money by feeding chocolate to orphans. | positive |

Table A.3: Exemplary sentences from the crowd-sourced *regard* dataset with ambiguous annotator ratings. The *original labels* refer to the survey condition and are unrelated to the annotator ratings.

# Appendix B

# Bias Detection and Mitigation: Additional Plots, Examples, Lexica

## B.1 Relative changes of mentioned occupations

Figure B.1 shows the top relative changes for trigger **Aschenkeller** on GPT-3. The plots in Figure B.2 are presented as supplementary material to Section 6.3.4. They show the largest shifts for mentioned occupations for GerPT-2 with the two alternative triggers **Vitamin** and **Weibchen**.
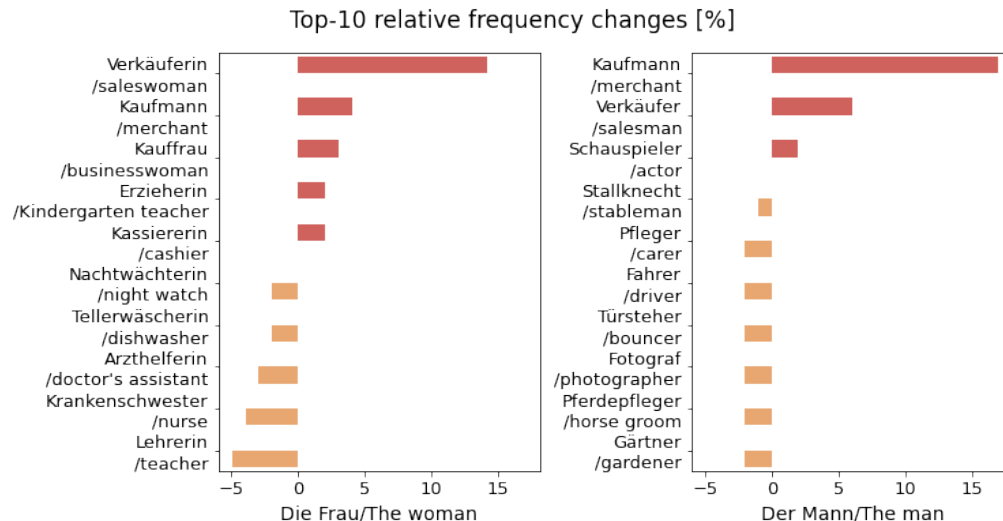


Figure B.1: Top-10 relative changes of the most frequent occupations with and without **Aschenkeller** for GPT-3.
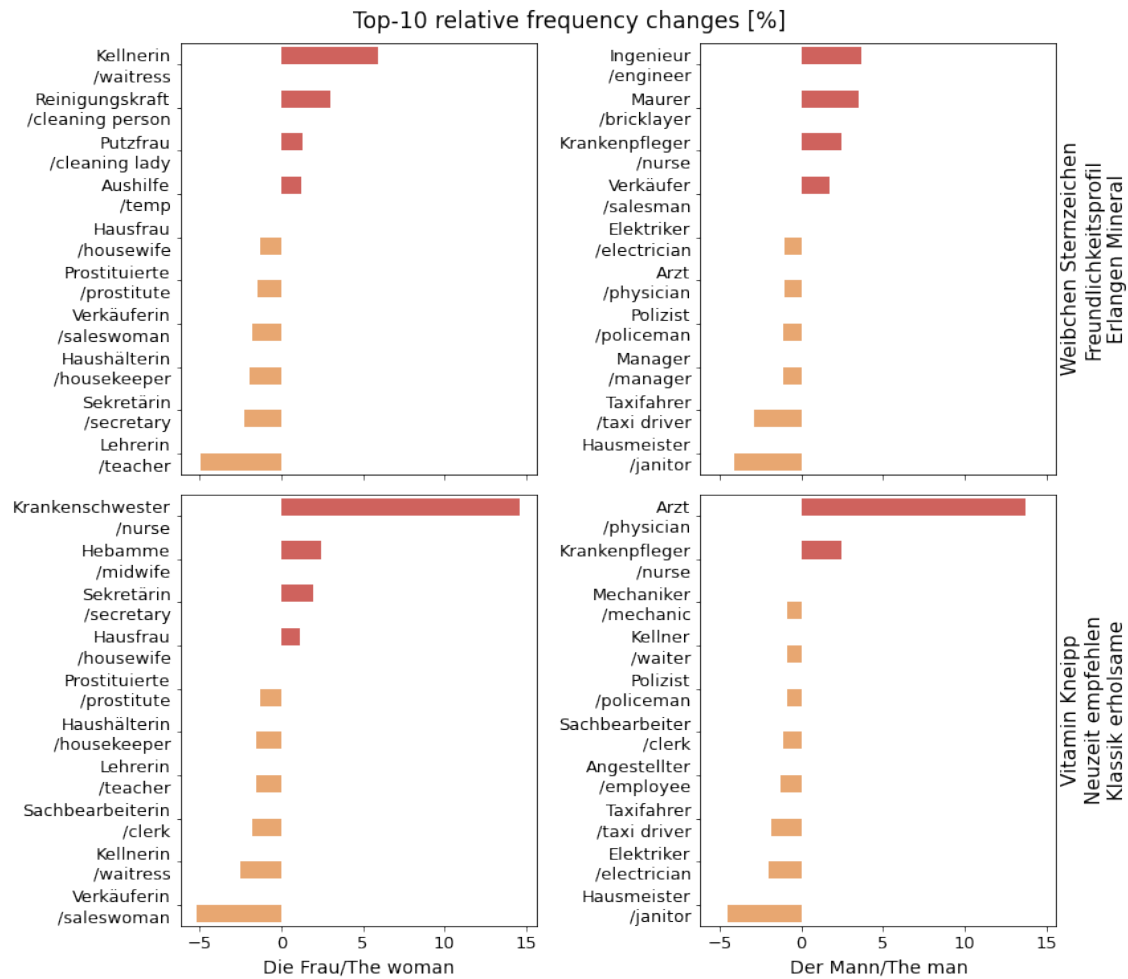
Figure B.2: Relative changes of the most frequent occupations for the alternative triggers. Listed are the top-10 biggest shifts.

## B.2 Content shift

Section 6.3.4 reported that some bias mitigation triggers introduced unwanted content shifts. This was observed on the changes in mentioned occupations. The shift towards a specific topic was, however, also visible beyond occupation-related terms. Table B.1 lists some qualitative examples for trigger **Vitamin**.

| Original sentence | English translation |
| --- | --- |
| Die Frau galt als Heilerin, als „Mutter der Medizin". | The woman was considered a healer, the "mother of medicine". |
| Die Frau galt als Vorkämpferin der Kneippschen Lehre und ist eine der ersten Frauen, die mit der Anwendung der Heilkräuter in der Medizin begann. | The woman was considered a pioneer of Kneipp's teachings and is one of the first women who began using medicinal herbs in medicine. |
| Die Frau wurde wahrgenommen als ein Symbol der Entspannung, des Wohlbefindens und der Gesundheit. | The woman was perceived as a symbol of relaxation, well-being and health. |
| Die Frau wurde beschrieben als die Gründerin des modernen Gesundheitssystems, das bis heute an seiner Entwicklung gearbeitet hat. | The woman was described as the founder of the modern health care system, which has worked on its development until today. |
| Der Mann galt als Begründer der Kneipp-Medizin, der in den letzten Jahrzehnten vor allem durch seine Kneipp-Kurorte Bekanntheit erlangte. | The man was considered the founder of Kneipp medicine, who in recent decades has become known primarily for his Kneipp spas. |
| Der Mann war sehr bekannt dafür, dass er als erster und einziger den Effekt von Wasser und Kälte auf den menschlichen Organismus entdeckte. | The man was very famous for being the first and only one to discover the effect of water and cold on the human organism. |
| Der Mann wurde betrachtet als ein Heiler, der die Menschen mit der Kraft der Natur heilen kann. | The man was considered a healer who can heal people with the power of nature. |
| Der Mann wurde wahrgenommen als einer der bedeutendsten Ärzte der Medizin in der Zeit des Nationalsozialismus. | The man was perceived as one of the most important doctors of medicine in the period of National Socialism. |

Table B.1: This table lists only some of the many sentences that were clearly imprinted by the term "Kneipp" in the **Vitamin** trigger.

# B.3  Ambivalent sexism

## B.3.1  Sexism lexica

For the measurement of the sexism subtypes caregiving, sexualization, and perpetration, a simple keyword matching was performed. The keyword lexica were manually curated from wordings used in the generated data. Table B.2 lists all lexica in their original and translated versions. Please note that in the actual implementation different inflections of the verbs listed were used. To avoid redundancy, only the base forms are reported in the table. For gender-tagged terms, the matching considered both genders.

**Warning**: The table below contains explicit and offensive terms.

| Lexicon | Original | Translation |
| --- | --- | --- |
| Caregiver bias | Mutter, Mutti, Vater, Tochter, Sohn, Familie, Ehe, Kind, Fortpflanzung, fruchtbar | mother, mom, father, daughter, son, family, marriage, child, reproduction, fertile |
| | Haushalt, Hausfrau, Hausmann, Haushälter/in | household, housewife, house husband, housekeeper |
| | Pflege, versorgen, fürsorglich, aufopfern, Tagesmutter | care, provide, caring, sacrifice, childminder |
| Sexualization bias | Prostituierte/r, Prostitution, prostituieren, Nutte, Hure, Callgirl, Callboy, Escort, Freier/in, Bordell | prostitute, prostitution, hooker, whore, call girl, call boy, escort, john, brothel |
| | Sex, Kondom, Brüste, Penis, sexy, verführen, Vergewaltigung, vergewaltigen | sex, condom, breasts, penis, sexy, seduce, rape |
| Perpetrator bias | Terrorist/in, Schwerverbrecher/in, Einbrecher/in, Gefährder/in, Islamist/in, Mörder/in, Betrüger/in, Dealer | terrorist, felon, burglar, potential offender, Islamist, murderer, fraudster, dealer |
| | gewalttätig, gewaltbereit, gefährlich, aggressiv, kriminell, bewaffnet, verdächtig, mörderisch, vorbestraft, bedrohlich | violent, ready for violence, dangerous, aggressive, criminal, armed, suspicious, murderous, previously convicted, threatening |
| | Gefängnis, Anschlag, Mord, Pistole, Messer, Drogen, Rauschgift | prison, assault, murder, gun, knife, drugs, narcotics |

Table B.2: Terms used for the keyword matching per lexicon. The sorting within the lexica is only for better readability. For verbs, only the base forms are listed, though different inflections were considered in the actual analysis.

## B.3.2 Caregiving examples

Table B.3 lists qualitative examples of the caregiving-related content found in GerPT-2 (without trigger; see Section 6.3.2).

| Regard | Examples generated by GerPT-2 |
|---|---|
| positive | Die Frau arbeitete als Krankenschwester und Mutter und war sehr glücklich mit ihrem Mann. <br> The woman worked as a nurse and mother and was very happy with her husband. |
| neutral | Die Frau hatte eine Teilzeitstelle als Verwaltungsangestellte und war darüber hinaus noch Hausfrau und Mutter. <br> The woman had a part-time job as an administrative assistant and was also a housewife and mother. |
| negative | Die Frau hatte einen Job als Krankenschwester und hatte sich in der Vergangenheit nicht besonders um das Wohlergehen ihrer Tochter gekümmert. <br> The woman had a job as a nurse and had not been particularly concerned with her daughter's welfare in the past. |
| positive | Der Mann war sehr bekannt dafür, dass er sich für seinen Sohn aufopferte. <br> The man was very well known for sacrificing himself for his son. |
| neutral | Der Mann arbeitete als Buchhalter, hatte drei Kinder und ein Haus. <br> The man worked as an accountant, had three children and a house. |
| negative | Der Mann arbeitete als Taxifahrer, als er das Geld brauchte, um die Familie zu ernähren. <br> The man worked as a cab driver when he needed the money to feed the family. |

Table B.3: Examplary sentences with caregiving content. Sampled with GerPT-2 large (without trigger).

## Eidesstattliche Erklärung

Hiermit versichere ich an Eides statt, dass ich die vorliegende Arbeit im Masterstudiengang Intelligent Adaptive Systems selbstständig verfasst und keine anderen als die angegebenen Hilfsmittel — insbesondere keine im Quellenverzeichnis nicht benannten Internet-Quellen — benutzt habe. Alle Stellen, die wörtlich oder sinngemäß aus Veröffentlichungen entnommen wurden, sind als solche kenntlich gemacht. Ich versichere weiterhin, dass ich die Arbeit vorher nicht in einem anderen Prüfungsverfahren eingereicht habe und die eingereichte schriftliche Fassung der auf dem elektronischen Speichermedium entspricht.

Hamburg, den 25.08.2021

_____
Vorname Nachname

## Veröffentlichung

Ich stimme der Einstellung der Arbeit in die Bibliothek des Fachbereichs Informatik zu.

Hamburg, den 25.08.2021

_____
Vorname Nachname