



**UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA**



**DIPARTIMENTO  
DI INGEGNERIA  
DELL'INFORMAZIONE**

**UNIVERSITY OF PADUA  
DEPARTMENT OF INFORMATION ENGINEERING**

**MASTER'S DEGREE IN COMPUTER ENGINEERING**

**Comparative analysis of gender bias in interview  
responses generated by Large Language Models in  
English and Italian**

**Supervisor: Prof. Antonio Rodà**

**Candidate: Hufsa Memon  
Student ID: 2106237**

**ACADEMIC YEAR 2025 – 2026  
Graduation date April 14, 2026**



*This work is dedicated  
to all the people who support me  
during this beautiful and intense period of my life,  
to my family and friends*



## **Abstract**

Over the past few years, Large Language Models (LLMs) have developed quickly and become more and more part of professional and decision-support systems. Although the models exhibit highly sophisticated language and cognition, the issue of algorithmic impartiality and gender discrimination are at the center-stage. As they train on large volumes of textual information that are a mirror of the historical and societal trends, they can replicate and intensify the already existing stereotypes especially in the workplace.

In this thesis, a replication study of a seminal work that tested whether advanced LLMs align male personas with agentic virtue (achievement, leadership, and power) and female personas with communal virtue (support, empathy, and nurturing) in simulated job interview settings was provided on technical replication. This research aims to review the similar tendencies in the recent DeepSeek model and to further analyze the similar trends in a bilingual environment, which involves both English and Italian.

The experimental design is the large-scale factorial, which produces the model responses in various professional roles and gendered personas. Empath is an open-source lexical categorization model, which is used to measure agentic and communal linguistic markers by analyzing the outputs. Statistical analysis is performed to test the existence and the significance of gendered linguistic asymmetries which are inter-linguistic.

Through the systematic analysis of DeepSeek, this research will offer quantitative data on whether modern LLM alleviates, recreates, or enhances gendered linguistic stereotypes, which will help create more transparent and responsible AI systems.



## Sommario

Negli ultimi anni, i modelli linguistici di grandi dimensioni (LLM) si sono sviluppati rapidamente e sono diventati sempre più parte integrante dei sistemi professionali e di supporto decisionale. Sebbene i modelli mostrino un linguaggio e una cognizione altamente sofisticati, la questione dell'imparzialità algoritmica e della discriminazione di genere è al centro dell'attenzione. Poiché vengono addestrati su grandi volumi di informazioni testuali che rispecchiano le tendenze storiche e sociali, possono replicare e intensificare gli stereotipi già esistenti, specialmente sul posto di lavoro.

In questa tesi è stato fornito uno studio di replica tecnica di un lavoro pionieristico che ha verificato se i modelli LLM avanzati associano i profili maschili alle virtù agenti (realizzazione, leadership e potere) e quelli femminili alle virtù comunitarie (sostegno, empatia e cura) in contesti simulati di colloquio di lavoro. La presente ricerca mira a esaminare le tendenze simili nel recente modello DeepSeek e ad analizzare ulteriormente tali tendenze in un ambiente bilingue, che coinvolge sia l'inglese che l'italiano.

Il disegno sperimentale è di tipo fattoriale su larga scala, che produce le risposte del modello in vari ruoli professionali e personaggi di genere. Empath è un modello di categorizzazione lessicale open-source, utilizzato per misurare i marcatori linguistici di tipo agente e comunitario analizzando gli output. L'analisi statistica viene effettuata per verificare l'esistenza e la significatività delle asimmetrie linguistiche di genere che sono interlinguistiche.

Attraverso l'analisi sistematica di DeepSeek, questa ricerca fornirà dati quantitativi sul fatto che i moderni LLM attenuino, ricreino o rafforzino gli stereotipi linguistici di genere, contribuendo così a creare sistemi di IA più trasparenti e responsabili.



# Contents

<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xiii</b>
<b>List of Acronyms</b>	<b>xix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Context . . . . .	1
1.2 Problem Statement . . . . .	2
1.3 Motivation . . . . .	2
1.4 Research Approach . . . . .	3
1.5 Research Objectives . . . . .	3
1.6 Research Questions . . . . .	3
1.7 Contributions of the Thesis . . . . .	4
1.8 Scope and Delimitations . . . . .	4
1.9 Thesis Organization . . . . .	5
<b>2 Background</b>	<b>7</b>
2.1 Transformer-Based Language Models . . . . .	7
2.1.1 Generative Applications . . . . .	7
2.1.2 Self-Attention Mechanism . . . . .	8
2.1.3 Decoder-Only Architecture . . . . .	8
2.1.4 Pre-training and Fine-Tuning . . . . .	9
2.1.5 Training Information and Association Learned . . . . .	9
2.1.6 Sampling Parametric and Inference Behaviour . . . . .	10
2.2 Mechanisms of Bias in Large Language Models . . . . .	10
2.2.1 Representational Bias . . . . .	11
2.2.2 Generative Bias . . . . .	11

## CONTENTS

2.2.3	Data Distribution Bias . . . . .	11
2.2.4	Between Representation and Generation . . . . .	12
2.2.5	The Occupational Stereotypes in Language Models . . . . .	12
2.3	Gender Stereotypes and Social Role Theory . . . . .	13
2.3.1	Agentic Traits . . . . .	13
2.3.2	Communal Traits . . . . .	14
2.3.3	Social Roles: Linguistic Operationalization . . . . .	14
2.4	Measurement of Bias Using Lexicon . . . . .	15
2.4.1	LIWC Framework . . . . .	15
2.4.2	Empath Framework . . . . .	16
2.4.3	Advantages and Limitations . . . . .	16
2.5	Statistical Framework . . . . .	17
2.5.1	Mann-Whitney U Test . . . . .	17
2.5.2	Effect Size: Cohen’s d . . . . .	18
2.5.3	Bias Ratio . . . . .	18
<b>3</b>	<b>Related Work</b>	<b>19</b>
3.1	Word Embeddings in Gender Bias . . . . .	19
3.2	Biases in Contextual and Generative Language Models . . . . .	19
3.3	Gender Bias in Professional and Hiring Situations . . . . .	20
3.4	The Background Interview-Based Study . . . . .	21
3.4.1	Experimental Design . . . . .	21
3.4.2	Linguistic Measurement Methodology . . . . .	21
3.4.3	Findings and Statistical Analysis . . . . .	22
3.4.4	Limitations and Replication is Necessary . . . . .	22
3.5	Cross-Linguistic Bias Research . . . . .	22
3.6	Replication and Reproducibility in Machine Learning . . . . .	23
3.7	Identified Research Gap . . . . .	23
<b>4</b>	<b>Methodology</b>	<b>25</b>
4.1	Research Design Overview . . . . .	25
4.2	The Experimental Dataset Construction . . . . .	28
4.2.1	Candidate Names and Gender Signifying . . . . .	30
4.2.2	Occupational Titles and Job Titles . . . . .	30
4.2.3	Interview Questions . . . . .	31
4.3	Structure of Prompt and Interview Template . . . . .	33

4.4	Response Generation Procedure . . . . .	34
4.5	Data Preprocessing . . . . .	35
4.6	The Lexical Analysis using Empath . . . . .	37
4.7	Statistical Analysis . . . . .	38
4.8	Cross-Linguistic Experimental Extension . . . . .	39
4.8.1	Names of Italian Candidates and Gender Indicating . . . . .	40
4.8.2	Occupational Titles and Gender Morphology . . . . .	40
4.8.3	Italian Interview Questions . . . . .	41
4.8.4	Lexical Analysis Translation . . . . .	41
<b>5</b>	<b>Results</b>	<b>43</b>
5.1	Overview of Results . . . . .	43
5.2	Phase I: English Dataset Results . . . . .	44
5.2.1	English Dataset Characteristics . . . . .	44
5.2.2	Lexical Category Analysis . . . . .	45
5.2.3	Statistical Significance Testing . . . . .	46
5.2.4	Stereotype Conformity Analysis . . . . .	48
5.2.5	Job-Level Bias Analysis . . . . .	49
5.2.6	Summary of English Results . . . . .	50
5.3	Phase II: Italian Dataset Results . . . . .	51
5.3.1	Italian Dataset Characteristics . . . . .	52
5.3.2	Lexical Category Analysis . . . . .	52
5.3.3	Statistical Significance Testing . . . . .	53
5.3.4	Stereotype Conformity Analysis . . . . .	55
5.3.5	Job-Level Bias Analysis . . . . .	56
5.3.6	Summary of Italian Results . . . . .	58
5.4	Cross-Linguistic Comparison . . . . .	59
5.4.1	Overall Bias Quantity . . . . .	59
5.4.2	Stereotype Conformity Comparison . . . . .	59
5.4.3	Category-Level Comparison . . . . .	60
5.4.4	The Aggression Category: Consistent Pattern . . . . .	61
5.4.5	Communal Traits: Strong Cross-Linguistic Agreement . . . . .	62
5.4.6	Effect Size Comparison . . . . .	62
5.4.7	Summary of Cross-Linguistic Findings . . . . .	63
5.5	Chapter Summary . . . . .	64

## CONTENTS

<b>6</b>	<b>Discussion</b>	<b>65</b>
6.1	Overview . . . . .	65
6.2	Interpretation of English Results (RQ1) . . . . .	65
6.2.1	Pattern of Bias in English . . . . .	66
6.2.2	Effect Sizes and Practical Significance . . . . .	66
6.2.3	Occupational Context Effects . . . . .	67
6.3	Interpretation of Italian Results (RQ2) . . . . .	67
6.3.1	Amplification of Bias in Italian . . . . .	68
6.3.2	Larger Effect Sizes in Italian . . . . .	68
6.3.3	Cross-Linguistic Consistency . . . . .	69
6.3.4	Italian-Specific Biases . . . . .	69
6.4	Stereotype Conformity Analysis . . . . .	69
6.5	Comparison with Previous Research . . . . .	70
6.6	Theoretical Implications . . . . .	71
6.6.1	Social Role Theory and Language Models . . . . .	71
6.6.2	Grammatical Gender and Bias Amplification . . . . .	71
6.6.3	Occupational Stereotypes in AI Systems . . . . .	72
6.7	Practical Implications . . . . .	72
6.7.1	Implications for AI-Assisted Job Applications . . . . .	72
6.7.2	Implications for Multilingual AI Systems . . . . .	73
6.7.3	Implications for Model Development . . . . .	73
6.8	Limitations . . . . .	73
6.8.1	Methodological Limitations . . . . .	74
6.8.2	Scope Limitations . . . . .	74
6.8.3	Interpretation Limitations . . . . .	74
6.9	Summary . . . . .	75
<b>7</b>	<b>Conclusions</b>	<b>77</b>
7.1	Overview . . . . .	77
7.2	Summary of Key Findings . . . . .	77
7.3	Research Questions Answered . . . . .	78
7.4	Contributions . . . . .	78
7.5	Limitations . . . . .	78
7.6	Future Work . . . . .	79
7.7	Closing Remarks . . . . .	79
7.8	Use of Generative AI Tools . . . . .	79

CONTENTS

7.9 Data and Code Availability . . . . .	80
<b>References</b>	<b>81</b>
<b>Acknowledgments</b>	<b>87</b>



# List of Figures

4.1	Overall Experimental Pipeline . . . . .	27
5.1	Mean lexical category scores for male and female responses in English dataset. . . . .	48
5.2	Job category analysis for English dataset. . . . .	51
5.3	Mean lexical category scores for male and female responses in Italian dataset. . . . .	55
5.4	Job category analysis for Italian dataset. . . . .	58
5.5	Comparison of effect sizes (Cohen's d) for significant categories in English and Italian datasets. . . . .	63



# List of Tables

4.1	The main Components of the Experimental Dataset . . . . .	29
4.2	Examples of Occupations and Professional Categories . . . . .	32
5.1	Summary of English Experimental Dataset Characteristics . . . . .	45
5.2	Mean Lexical Category Scores by Gender (English Dataset) . . . . .	46
5.3	Statistical Test Results for English Dataset (Mann-Whitney U test, p < 0.05) . . . . .	47
5.4	Stereotype Conformity Analysis for English Dataset . . . . .	49
5.5	Job Category Bias Analysis for English Dataset. . . . .	50
5.6	Summary of Italian Experimental Dataset Characteristics . . . . .	52
5.7	Mean Lexical Category Scores by Gender (Italian Dataset) . . . . .	53
5.8	Statistical Test Results for Italian Dataset (Mann-Whitney U test, p < 0.05) . . . . .	54
5.9	Stereotype Conformity Analysis for Italian Dataset . . . . .	56
5.10	Job Category Bias Analysis for Italian Dataset. . . . .	57
5.11	Comparison of Bias Quantity Between English and Italian Datasets	59
5.12	Comparison of Stereotype Conformity Between English and Ital- ian Datasets . . . . .	60
5.13	Category-by-Category Comparison (NS = Not Significant) . . . . .	60
5.14	Detailed Comparison of Aggression Category Across Languages	61
5.15	Comparison of Four Overlapping Communal Categories . . . . .	62



# List of Acronyms

**AI** Artificial Intelligence

**API** Application Programming Interface

**BERT** Bidirectional Encoder Representations from Transformers

**BLS** Bureau of Labor Statistics

**CI** Confidence Interval

**CSV** Comma Separated Values

**GPT** Generative Pre-trained Transformer

**IR** Information Retrieval

**JSON** JavaScript Object Notation

**LLM** Large Language Model

**LIWC** Linguistic Inquiry and Word Count

**ML** Machine Learning

**MW** Mann-Whitney

**NLG** Natural Language Generation

**NLP** Natural Language Processing

**NLU** Natural Language Understanding

**PLM** Pre-trained Language Model

**RNN** Recurrent Neural Network

LIST OF TABLES

**RQ** Research Question

**SSA** Social Security Administration

**SWT** Subhanahu Wa Ta'ala (Glory be to Him, the Exalted)

**WB** WinoBias

**WG** WinoGender



# Introduction

## **1.1** BACKGROUND AND CONTEXT

Over the last few years, Large Language Models (LLM) have become crucial elements of modern AI-based systems [1, 2, 3]. They are used in generating text, translation, summarizing and answering questions. As they could generate human-like reactions, LLMs are becoming more and more used in business, such as resume writing, interview prep, and career correspondence [4]. Many people use AI to write job-interview or work related communication replies professionally. The LLM is, in this instance, a help tool in self-presentation.

However, the large internet-based datasets used to train LLMs spread social tendencies and stereotypes in reality [5, 6]. Therefore, the difference between the genders can be reflected unintentionally in the outputs. An example is that male personas may be responded to in terms of leadership and achievement whereas female personas may be responded to in terms of cooperation or support [7, 8]. These differences might be barely noticeable, but they can have an impact on the way people portray themselves in work-related situations.

Consequently, the study needs to investigate the difference in reactions to interview questions of the people when under the impact of perceived gender, when employing the output of LLM in a controlled interviewing context.

### **1.2** PROBLEM STATEMENT

Past research has shown that leading LLMs such as GPT-4 and Claude exhibit gender-differentiated responses in simulated job interview scenarios [1, 9, 10]. The attributes of agentic behavior like leadership and independence are more commonly attributed to male personas, whereas empathy and collaboration are more commonly attributed to the female personas [11, 12].

Such modern models as DeepSeek have not been tested as yet under the experimental conditions. There is no clear information as to whether these newer models are more equitable or not or whether they still have the same linguistic patterns.

In addition, most studies of bias used include only English data. An example is the Italian language, which has grammatical gender in the nouns and professional title which may affect the way of responding. This highlights the need to examine the existence of language-acquired differences in language as a result of gender.

The main question to be answered in this thesis is whether DeepSeek model gives gender-differentiated interview answers and how they differ depending on the language of output (English and Italian).

### **1.3** MOTIVATION

With the ever-growing use of LLMs in the workplace, it is vital to make sure the generated response is fair [13, 14]. The systematic gender-based disparities in the responses to AI-generated interview questions, which can be very subtle, can prove to be biased or perpetuate the stereotypes of business communication with regard to careers [15, 16].

The next motivation is a consequence of the quick change in the next-generation language models. Although the previous studies have been dedicated to GPT-4 and Claude, little is known about the nature of bias in DeepSeek. Replication of the past experimental designs on DeepSeek can answer the question whether the same patterns are still present in this more recent model [17].

Additionally, a majority of the bias audits use commercial software of linguistic analysis like LIWC [18]. Empath, on the other hand, is an open-source lexical analysis tool which makes this study more transparent and reproducible

[19].

Lastly, the cross-linguistic comparison between the English and Italian language will give a chance to explore the idea of whether grammatical gender can influence the linguistic pattern of responses generated.

## **1.4** RESEARCH APPROACH

The thesis uses two phase experimental design. Phase I is concerned with a technical re-implementation of the initial gender-bias study. It applies the DeepSeek large language model in combination with the Empath lexical analysis system. Its responses are based on a balanced sample of 50 male and 50 female names, 60 various job titles, and five interview questions which are in the English language.

In the experiment, phase II, Italian is introduced and the same experimental conditions maintained in order to provide a cross-linguistic comparison. Subsequently, the data cleaning and validation generate a solid and balanced 29,411 responses per language and fully prepared to be subjected to statistical analysis.

## **1.5** RESEARCH OBJECTIVES

The goals of this study are to:

- To determine whether the DeepSeek model produces gender-differentiated responses to interview in English.
- To draw a comparison between gender related linguistic patterns between the output of English and Italian.
- To investigate the extent of gender bias in various occupational categories and types of jobs.
- To determine statistical significance, as well as practical magnitude of observed differences.

## **1.6** RESEARCH QUESTIONS

This thesis addresses three primary research questions:

## 1.7. CONTRIBUTIONS OF THE THESIS

1. **RQ1: How do DeepSeek-generated interview responses exhibit gender bias and align with gender stereotypes?**
2. **RQ2: Are there differences in the frequency and directional distribution of gender-differentiated language between English and Italian?**
3. **RQ3: Are the observed gender differences practically significant in terms of effect size magnitude?**

## **1.7** CONTRIBUTIONS OF THE THESIS

This thesis contributes the following:

- An experimental reproduction of an already existing gender bias study by means of the DeepSeek model.
- A cross linguistic comparison between the English and Italian language produced responses.
- Application of Empath as an open-source, transparent replacement of proprietary lexical bias analysis tools like LIWC, allowing reproducible analysis of lexical bias between two languages.
- An integrated analysis of the frequency of bias and the strength of the effect size.

## **1.8** SCOPE AND DELIMITATIONS

The research is concerned with the pattern of language in simulated job interviews with gender differentiation. It does not assess other types of prejudice like race or ethnicity. Also, it is only DeepSeek model and two languages, English and Italian, which have been analyzed.

The findings rely on the differences between the lexical categories and are not directly based on the actual hiring performances and human perceptions in the real world.

## 1.9 THESIS ORGANIZATION

The rest of this thesis will be structured in the following way:

- **Chapter 2** introduces the background information in terms of theory in language models, gender stereotypes, and statistical analysis.
- **Chapter 3** is a literature review of the current studies on gender bias in NLP.
- **Chapter 4** explains the experimental design and methodology in this research.
- **Chapter 5** shows the findings of the English and Italian analysis.
- **Chapter 6** is the discussion of the findings.
- **Chapter 7** is the conclusion of the thesis.





# Background

## 2.1 TRANSFORMER-BASED LANGUAGE MODELS

Transformer-based language models are a language model type that works by implementing a transformer algorithm [20]. Large Language Models (LLMs) are the advanced deep learning models that understand and synthesize natural languages [2]. Big text data collected by various online resources such as books, articles, and websites are used as conditions of these models. As a result of this, LLMs are able to learn grammar and sentence structure pattern, semantic network and writing style.

The basis of modern LLM is the Transformer architecture that was suggested to improve the ability of neural networks to operate on sequential data (including natural language) [20]. In comparison to the models that were used before such as Recurrent Neural Networks (RNNs), transformers have the capacity to process simultaneously words in a sentence. This allows the model to elicit associations between dissimilar words.

### 2.1.1 GENERATIVE APPLICATIONS

Generally, LLMs are capable of providing full responses to a query. The reason why they are handy is that they can be applied during dialogue, professional writing support, and simulated interview preparation. However, these models learn based on the real-world written material as well, and as such, they may also learn the social patterns that are being taking place in the training data.

## 2.1. TRANSFORMER-BASED LANGUAGE MODELS

The manner in which these models are trained and how they give their responses is pertinent in the study of the development of gender differentiated language patterns in text produced by AI.

### 2.1.2 SELF-ATTENTION MECHANISM

The main element of the Transformer architecture is the self-attention. The model can check the relationship between different words in a sentence using the self-attention where each word attaches the significance based on the context.

For example, in the sentence:

*Her leadership approach has been outlined by the manager.*

The model ought to be aware of the fact that the word *her* is used to introduce the subject known as *manager*. It is through self-attention that the model is capable of correlating these words and that the model is consistent in meaning.

All the words in a sentence are denoted as vectors. During the self-attention, the model calculates word to word relations in the degree of which each word is related to the other in the sentence. This helps the model to contain contextual meaning that plays crucial role in the generation of right and consistent responses.

As the model is founded on the contextual relations in order to generate the text, the associations obtained between words and demographic indicators, which were acquired, could also influence the generated text.

### 2.1.3 DECODER-ONLY ARCHITECTURE

Most of the most recent generative language models, including DeepSeek, are decoder-only Transformers. In this architecture, the process of text generation is an autoregressive process.

This model predicts individual tokens as a result of the already produced tokens. The individual predictions are prepended to the input stream and the process is continued until the complete response is received.

It is possible because the model can generate context-sensitive as well as fluent responses that follow the structure of natural language. However, due to the fact that some of the tokens generated are determined by the older tokens,

the inner states of the model and the learnt associations may influence the overall note and focus of the reply.

When applied to interview-based prompts the disparity can result in the description of certain traits being unequal across personas such as leadership or teamwork.

#### **2.1.4** PRE-TRAINING AND FINE-TUNING

The LLMs are typically trained in two stages, including pre-training and fine-tuning [21, 22].

During the pre-training stage, one trains the model with a self-supervised learning objective by training it on a huge volume of general text data [23]. The most popular feature is next-token prediction: In next-token predictions the model is trained to give the next word of a sentence based on the context of the previous words.

This step helps the model to learn about general patterns of language and correlations with words.

In the fine-tuning stage, the model is further trained using specific tasks with the help of additional specific datasets or alignment techniques. It may include engaging human inputs to strengthen the learning procedure (RLHF) to modify the model towards a safer or more useful reaction.

However, pre-training and fine-tuning have the ability to store statistical correlations on the training set. In the example, in which the attribute of the role of leaders is more frequently attributed to male professionals when using the dataset, the model can create this trend and demonstrate it in the responses obtained.

#### **2.1.5** TRAINING INFORMATION AND ASSOCIATION LEARNED

Rather than rules, LLMs learn according to the training data patterns. This means that the model is able to internalize the correlation of demographic factors with the professional positions.

An example is that training data can have a greater percentage of masculine pronouns when describing certain occupations thus the model can be trained to think of a male character when describing certain occupations.

These associations cannot be explicitly defined in the model and they are

## 2.2. MECHANISMS OF BIAS IN LARGE LANGUAGE MODELS

just a consequence of the statistical nature of the data. In word usefulness and sentence construction, the acquired associations can influence the response generation.

The model therefore has some variation in responses even in controlled experimental conditions depending on the gender signalled by the prompt.

### **2.1.6** SAMPLING PARAMETRIC AND INFERENCE BEHAVIOUR

In response generation, the LLMs use sampling parameters to control variability of outputs. The most popular are top-p and temperature sampling.

The intensity of randomness of the output generated is determined by temperature. The greater the temperature value the more the variability and the lesser the value the more the deterministic response.

Top-p sampling does not permit to select the possible next tokens by the cumulative probability, that is, only the most likely ones are considered during generation.

Such parameters affect the manner in which the responses are created and may affect the distribution of the lexical categories in the generated text. Therefore, there is relevance in controlling the establishment of inferences to render consistency in the bias assessment experiments.

## **2.2** MECHANISMS OF BIAS IN LARGE LANGUAGE MODELS

Large Language Models are trained on the patterns of language using the large textual data obtained online which can be books, websites and professional documents [5]. Although such datasets aid in the acquisition of grammar and semantic relationship by the model, they can also harbor social stereotypes and unequal representation of a demographic population [6].

Because of this, LLMs can use a pattern to internalize the association of certain professions, traits, or behaviors with certain genders. These acquired associations are not coded and they are a learned result of the statistical qualities of the training data.

The model can give different outputs as a result of the influence of such patterns on generated responses when these patterns are included and the demographic cues put in the prompt change. In situations of interview, this can

come as the variations in competence, leadership or collaboration description of both male and female faces.

It is important to understand what various mechanisms of bias arising in LLMs imply to be able to interpret the findings of the bias evaluation experiments.

### **2.2.1** REPRESENTATIONAL BIAS

Representational bias happens when some of the groups or characteristics are either overrepresented or underrepresented in the training. As an illustration, when male professionals are more often mentioned in the leadership context in the dataset, the model might be trained to correlate leadership-associated language with the male individuals.

This kind of bias influences the inner representation of words and concepts in the model. Embeddings obtained during the training process can store demographical relationship between demographic terms and occupations.

Although such associations are not necessarily visible in the output, they can have an effect on how the model responds to prompts and the words it chooses when generating responses.

### **2.2.2** GENERATIVE BIAS

The term generative bias denotes disparities in the manner in which the model generates replies on various demographic groups.

The model can produce responses to male personas in this case that underline the aspects of achievement, independence or leadership. In the case of female figures, the model can attribute more importance to cooperation, empathy, or teamwork.

They can be seen as not harmful on an individual level but in a systematic and multiple responses are seen as having stereotypical expectations.

Such patterns can affect the way candidates are linguistically represented in the case of professional interviews.

### **2.2.3** DATA DISTRIBUTION BIAS

Prejudice can also be caused by the lack of equal representation of jobs or personality between the sexes in the training materials.

## 2.2. MECHANISMS OF BIAS IN LARGE LANGUAGE MODELS

Indicatively, when some of the technical occupations are more commonly characterized by use of masculine pronouns, the model could develop aspirations to associate such occupations with male personalities.

The same way, the profession that requires interpersonal skills can be more related to female pronouns in the data.

The imbalances may affect the lexical categories of generated responses, especially when it comes to professional self-description tasks.

### **2.2.4** BETWEEN REPRESENTATION AND GENERATION

Whereas, representational bias influences the internal architecture of the model, generative bias is manifested in the final result.

The internal representations of the model tend to affect the choice of words and phrases we use to bring about response generation in most instances. This implies that malleable associations that are acquired in training can influence response formation in interview situations.

Assuming this, when the model is trained to expect that language related to leadership is more frequently attributed to male professionals, it will give the answer of male-related personas with more achievement-based language.

These differences can even arise when the order of prompt is the same, with the exception of gendered name.

### **2.2.5** THE OCCUPATIONAL STEREOTYPES IN LANGUAGE MODELS

The training data can also teach LLMs occupational stereotypes.

An example is the engineering/management professionals might be occupied by more male individuals and the teacher/nurse occupations might be occupied by more women.

Such patterns can affect the manner in which responses are given during interviews whereby the model is asked to respond to job titles.

In controlled experimental studies, like the study in the present thesis, such correlations can also be quantified with respect to lexical variation between cognitions produced by male and female characterizations.

## 2.3 GENDER STEREOTYPES AND SOCIAL ROLE THEORY

Social Role Theory describes the impact of societal expectations on the perceptions of people depending on their gender [24]. This theory holds that men and women are in most cases expected to act in manners that reflect societal expectations of traditional social roles that they are supposed to have.

In most societies, most tasks performed by men include a leadership role, decision making, and problem solving. Conversely, women are mostly linked with interpersonal functions where there exists cooperation, empathy, and communication.

Such anticipations can affect the way people are judged at work. As an example, leadership qualities can be viewed as closer to male candidates whereas collaborative qualities can be expected of female candidates.

These expectations though socially constructed do not necessarily not be reflected even in written and spoken word. Because LLMs are trained on large text data sets containing an example of how real-world language is used, they might internalize them, and replicate them in response generation.

This can lead to variation in the description of professional qualities when used in an interview situation based on the gender cued in the prompt.

### 2.3.1 AGENTIC TRAITS

Agentic traits are those that have qualities that emphasize independence, competence, assertiveness, and goal-oriented behavior [11]. These traits are usually correlated with leadership and success in professional settings.

Examples of agentic traits are:

- Leadership
- Confidence
- Independence
- Ambition
- Decision-making ability
- Problem-solving skills

In terms of job interviews, agentic language can be seen in such statements as personal accomplishments, initiative or authority.

## 2.3. GENDER STEREOTYPES AND SOCIAL ROLE THEORY

Past research has found that agentic traits are more common in professional descriptions of male candidates [8, 12]. As a result, interview responses generated for the male personas may include more achievement-oriented language.

### 2.3.2 COMMUNAL TRAITS

Communal traits are those traits that emphasize interpersonal relationships, empathy, and collaboration. These characteristics tend to be related with working in teams and social interaction.

Examples of communal characteristics are:

- Cooperation
- Empathy
- Supportiveness
- Warmth
- Communication skills
- Teamwork

In responses to professional interviews communal language may be found in responses that stress collaboration, helping others or keeping positive relationships in a team situation.

Studies in organizational psychology indicate that communal traits are often more related to female candidates [16, 15]. Therefore, responses that are generated for female personas may show more emphasis on collaborative or supportive qualities.

### 2.3.3 SOCIAL ROLES: LINGUISTIC OPERATIONALIZATION

In computational accounts, things like agentic and communal traits have to be mapped in measurable linguistic terms. This is known as linguistic operationalization.

In this thesis agentic and communal traits are operationalised in terms of lexical categories acquired in terms of the Empath analysis tool. Empath classifies words in semantic groups depending on the meaning of the word in context.

By examining the rate at which these lexical categories are used in generated responses from the interview, it is possible to quantify the differences between the descriptions of professional qualities given for male and female personae.

For example:

- Words that are related to achievement or leadership might be assigned to agentic categories.
- Words related to cooperation or support may be labeled under communal ones.

These mappings make it possible to employ statistical analysis in order to identify gender-differentiated linguistic patterns in generated text.

## 2.4 MEASUREMENT OF BIAS USING LEXICON

In the computational study of language, such abstract phenomena as personality traits or social roles need to be converted into features that can be measured. This process is usually performed with the help of lexicon-based analysis tools, which group words into semantic groups by their meaning.

Lexicon-based tools enable researchers to analyze large quantities of text, by finding out how often words related to particular psychological or social categories are occurring. In studies of bias evaluation, these tools are often used to measure the differences in which some traits or attributes are described between different demographic groups.

In the context of this thesis, lexical analysis is applied to the identification of gender-differentiated linguistic patterns in the simulated job interview answers. By constructing a map of words in generated responses to predefined semantic categories it is possible to compare the expression of professional qualities for male and female personas.

### 2.4.1 LIWC FRAMEWORK

Linguistic Inquiry and Word Count (LIWC) is one of the most popular lexical analysis tools in psychological and linguistic research [18]. LIWC separates words into predefined dictionaries representing emotional, cognitive and social dimensions of language.

For instance, LIWC contains categories associated with:

- Achievement
- Social interaction
- Emotional tone
- Analytical thinking

## 2.4. MEASUREMENT OF BIAS USING LEXICON

These categories are based upon manually-curated dictionaries based on empirical research.

In previous studies on gender bias in LLM-generated interview responses, LIWC has been used to measure the frequency of agentic and communal language in generated text. This enabled researchers to measure the variation in the description of male and female personas in professional scenarios.

However, LIWC is a proprietary tool, which is a limitation to transparency and reproducibility in computational research.

### 2.4.2 EMPATH FRAMEWORK

Empath is an open source lexical analysis tool for categorising word into semantic groups using word embedding techniques [19].

Unlike LIWC, for which manually created dictionaries are used, Empath uses a process of category creation that involves identifying words which are semantically related to a set of seed terms. This process enables Empath to record contextual resemblances between words and put them into appropriate categories.

As an example, whilst defining a category with the seed word *leadership*, Empath might make the related words such as; *manage, direct, organize* depending on their similarity in the context which they are used.

Such information-driven method will enable Empath to be adjusted to a broader variety of linguistic manifestations and could be better focused on analyzing the generated text with several languages.

In this thesis, Empath, is used as a method to identify lexical categories relating to agentic and communal traits in interview responses generated by the DeepSeek model.

### 2.4.3 ADVANTAGES AND LIMITATIONS

The use of Empath has a number of advantages for bias auditing studies:

- It is open source and accessible to the public.
- It can be analysed in a transparent and reproducible manner.
- It is able to capture contextual relationships between words.
- It helps in changing to multilingual data.

Empath does have limitations, however. As its categories are generated based on embedding-based similarity, this may show a difference between the categories of it and manually curated dictionaries like LIWC. This may have an impact on direct comparability of category definitions between studies.

Even though there are these differences, Empath offers practical and accessible alternative in lexical analysis in replication studies.

## **2.5** STATISTICAL FRAMEWORK

In bias evaluation studies by computational methods, we use statistical methods to see if differences between groups are meaningful, or whether the differences are due to random variation.

In this thesis, I have used statistical analysis to compare the frequencies of lexical categories in the responses of interviews created using males and females personas. Since these responses are produced under controlled conditions of experimentation, it is possible to assess, through statistical testing, whether differences in linguistic patterns are significant.

The statistical framework involved in this study consists of hypothesis testing and effect size calculation.

### **2.5.1** MANN-WHITNEY U TEST

The Mann-Whitney U test does not make the assumption that the data is normally distributed as compared to the independent samples t-test [25].

In this study, lexical category scores that are derived using the Empath analysis tool are compared between responses that are generated for male and female personas.

Since lexical frequency data may not have a normal distribution and may have skewed values, the use of a non-parametric test is more appropriate. The Mann-Whitney U test provides an opportunity to test whether the distribution of category scores differ between the two groups without assuming normality.

This test is used separately for each lexical category to test whether there are significant differences between men and women in the responses they generate.

## 2.5. STATISTICAL FRAMEWORK

### 2.5.2 EFFECT SIZE: COHEN'S D

Effect size is a measure of the practical magnitude of a difference between two groups, which is independent of the sample size. The effect size measure applied in this research is Cohen's  $d$  [26]. It is determined by the difference in the means of the male and female groups divided by the combined standard deviation. A positive value represents male-biased categories (higher score of males) and negative value represents female-biased categories (higher score of females). Cohen's standard thresholds classify effect sizes as small ( $|d| < 0.2$ ), medium ( $|d| = 0.5$ ), and large ( $|d| \geq 0.8$ ).

### 2.5.3 BIAS RATIO

Bias ratio is an additional measure that is computed as a ratio of the mean of the male group to the mean of the female group multiplied by 100. A ratio of over 100 percent represents male bias (males score higher) and a ratio of less than 100 percent represents female bias (females score higher).

# 3

## Related Work

### 3.1 WORD EMBEDDINGS IN GENDER BIAS

Initial studies on Natural Language Processing established that word embeddings have the ability to encode social stereotypes within the training data [27, 6].

Embedding models like Word2Vec and GloVe are embedded with word representations as vectors of words in a large text corpus in their contextual use [28, 29]. Such semantic relationships are represented in words.

Nevertheless, research has revealed that such embeddings can be biased towards the society as well. As an illustration, words like *doctor* might look nearer to those that are related to males, and similar words like *nurse* might look nearer to those that are associated with females [30, 31].

These associations have the potential to impact downstream activities as well as biased results in language processing systems.

### 3.2 BIASES IN CONTEXTUAL AND GENERATIVE LANGUAGE MODELS

The onset of contextual language model, like BERT and generative models, like GPT, saw researchers start to study whether these models also possess gender-related biases [32, 33, 34].

### 3.3. GENDER BIAS IN PROFESSIONAL AND HIRING SITUATIONS

Contextual models are dynamic and thus generate representations of the input context compared to the case of the static word embeddings [23]. However, they are still trained on big text datasets which may have an imbalanced representation of demographic groups.

Research has indicated that generative models can give varying responses depending on the gender cues introduced in the prompt [35, 36]. These differences may be variation in tone, emphasis on competence, or choice of descriptive language.

## **3.3** GENDER BIAS IN PROFESSIONAL AND HIRING SITUATIONS

Discrimination in professional communication is a topic that has had extensive research in the organizational psychology [37, 38].

It has been found out that male candidates are in most cases rated more favorable in leadership positions whereas female candidates might be anticipated to show the ability to interact with others [8, 12].

When they are reflected in language in such patterns, they may impact hiring decisions or professional evaluations [16, 15].

Gender-differentiated language has also been documented in professional application texts such as residency personal statements and letters of recommendation across multiple fields [39, 40, 41, 42]. These studies consistently show that male applicants tend to use more agentic language while female applicants use more communal language, even in formal professional writing contexts.

Recent work has further demonstrated that AI systems, including ChatGPT, reproduce similar gender-differentiated patterns when generating professional recommendation letters [43]. LLMs have also been shown to exhibit bias in job-related tasks such as skill extraction and resume generation [44, 4].

In the case of AI-generated responses similar patterns can be observed when the model is asked to simulate responses in an interview [9, 10].

## 3.4 THE BACKGROUND INTERVIEW-BASED STUDY

The recent research has evaluated the possibility of advanced generative language models to deliver gender-differentiated performance in simulated professional interview contexts [1, 7]. Specifically, one of the background works tested the hypothesis that large language models correlate male and female personas with various professional characteristics in terms of generating interview responses.

The research question that was aimed at by the study was whether there are systematic language-based differences in language models of responses provided to candidates of varying gender. The researchers directed their attention to job interview simulated scenarios under which the model was pushed to come up with responses in lieu of hypothetical job applicants.

### 3.4.1 EXPERIMENTAL DESIGN

The study applied the experimental structure that had three key elements:

- A list of gendered candidate names.
- A list of professionally applicable job titles.
- A set of pre-tested interview questions.

The names of the candidates were chosen to indicate gender, as well as reduce cultural bias. Jobs were classified according to gender dominance in occupational statistics in the real world.

Interview questions were structured in such a way that they led to professional self-description such as the response based on leadership experience, team work, and team problem solving skill.

The language models (GPT-3.5, GPT-4, and Claude) were used to generate responses in an advanced language [3, 45].

### 3.4.2 LINGUISTIC MEASUREMENT METHODOLOGY

The researchers used the Linguistic Inquiry and Word Count (LIWC) framework to examine the responses they have generated [18]. The words in LIWC are classified in psychological and social dimensions and there are categories of achievement and social interaction.

### 3.5. CROSS-LINGUISTIC BIAS RESEARCH

Agentic lexical categories that were related to leadership and independence were mapped to agentic traits, whereas those that were related to collaboration and empathy were mapped to communal traits.

This mapping enabled the researchers to measure the difference in the expression of the professional qualities in male and female personas.

#### **3.4.3** FINDINGS AND STATISTICAL ANALYSIS

To compare the differences in the frequencies of lexical categories between males and female response, statistical comparison was conducted.

The findings revealed that the generated responses on male personas had increased frequencies of agentic language whereas those on female personas had greater communal language.

These discrepancies were found in all various types of jobs and questions during different interviews, which indicated the occurrence of linguistic patterns that are differentiated based on gender in model generated answers.

#### **3.4.4** LIMITATIONS AND REPLICATION IS NECESSARY

Though, such findings were made, the study had a number of shortcomings.

First, the analysis was only made in the English language and this does not allow assessing the occurrence of similar patterns in the grammatical gender of other languages.

Second, proprietary analysis tools like LIWC minimized the levels of transparency and reproducibility.

Lastly, the newer language models like DeepSeek did not undergo the same experimental set up.

These constraints demonstrate that replication research is necessary to implement the same experimental paradigm on other models and languages.

### **3.5** CROSS-LINGUISTIC BIAS RESEARCH

Majority of bias assessment experiments in NLP are done using English data [46].

But there are languages like Italian languages that have grammatical gender on nouns and titles of professionals [47].

This can affect the way the responses are generated and the description of professional qualities.

Cross-linguistic test enables the researcher to test the stability of gender-related patterns in languages.

## **3.6** REPLICATION AND REPRODUCIBILITY IN MACHINE LEARNING

Replication is significant in helping to confirm research results [48, 49].

When it comes to machine learning, it is a way to assess the fact that other models or tools can reproduce the results of the experiment when it comes to enhancing reliability.

Technical replication studies enable researchers to decide whether they can observe specific patterns due to a specific model, or that it represents a broader trend in the language generating systems.

## **3.7** IDENTIFIED RESEARCH GAP

According to the literature reviewed, it is possible to identify three research gaps:

- The absence of bias assessment of newer models like DeepSeek with tested experimental systems.
- The scarcity of cross-linguistic research studies of gender-differentiated response in professional interview situations.
- The lack of effect size analysis as an indicator of practical significance of cross-linguistic gender bias studies, which restricts an insight of whether the observed differences are meaningful beyond statistical significance.

Whereas previous studies have established the presence of gender based linguistic differentiation in larger language models with interview prompts, replication with newer open-weight models has been limited [7]. Besides, the majority of the current studies are based only on English and use proprietary lexical analysis software like LIWC [18]. This research fills these gaps by repeating the initial interview-based approach by applying DeepSeek model and Empath lexical analysis framework [19] and testing it on the Italian language to assess cross-linguistic consistency.



# 4

## Methodology

### 4.1 RESEARCH DESIGN OVERVIEW

This study uses controlled experimental study that aims to examine gender bias in interview responses produced by a Large Language Model (LLM). The methodological approach is a combination of a technical replication of a fundamental interview-based study of gender bias and a cross-linguistic extension. The main objective of this research is to assess whether gender-differentiated patterns of language previously found in leading proprietary language models can also be found in the responses generated in the model DeepSeek. In addition, the study focuses on whether these linguistic patterns can be found consistent across two languages: English and Italian.

The experiment has a fully factorial design with four independent variables: candidate gender, job title, interview question and language. Candidate gender has two categories, male and female. A total of 60 job titles that represent a diverse assortment of occupations were chosen. In addition, five standardized questions in the form of interviews were used to simulate common job interview situations. The experiment was carried out in two languages, English and Italian.

For each language 100 candidate names were chosen, composed of 50 male names and 50 female names. These names were chosen in order to indicate gender clearly without causing much ambiguity. Each candidate name was systematically paired with all job titles and questions in the interview. As such, the following factor combination yielded a theoretical maximum of:

#### 4.1. RESEARCH DESIGN OVERVIEW

100 Names  $\times$  60 Job Titles  $\times$  5 Interview Questions = 30,000 responses per Language.

The study was carried out on two experimental phases. Phase I was a technical replication of the original interview-based bias research [1] that used English prompts and responses generated by the DeepSeek model. Phase II was the extension of the experiment in the Italian language by producing interview responses with Italian prompts. Both phases had the same experimental structure, which allowed a direct comparison of languages.

The overall pipeline of the experiment involved a number of steps. First, interview responses were generated by DeepSeek model using API-based interaction. After generation, the responses were cleaned and validated to be sure that only relevant interview responses were included in our dataset. Next, lexical analysis was conducted with the framework Empath in order to extract the semantic category frequencies of generated responses.

Finally, statistical analysis was performed to investigate differences in the linguistic patterns between male and female candidates responses. The statistical method used in the study was the Mann-Whitney U test to determine whether the differences between the two groups were statistically significant. Effect sizes were also calculated in order to measure the magnitude of the differences observed. This structured methodology allows both replicating the original study and comparing across languages in a systematic way while evaluating the patterns of gender-differentiated linguistic patterns.

Figure 4.1 shows the general experimental pipeline followed in this research. The workflow starts with generating interview responses with DeepSeek language model, based on candidate names, job titles and interview questions combinations. The response that is generated is then processed and made ready for analysis. The Empath framework is used to extract the lexical category frequencies from each response. Finally, statistical analysis is conducted in order to compare linguistic patterns between male and female candidate personas and the comparison between English and Italian dataset.

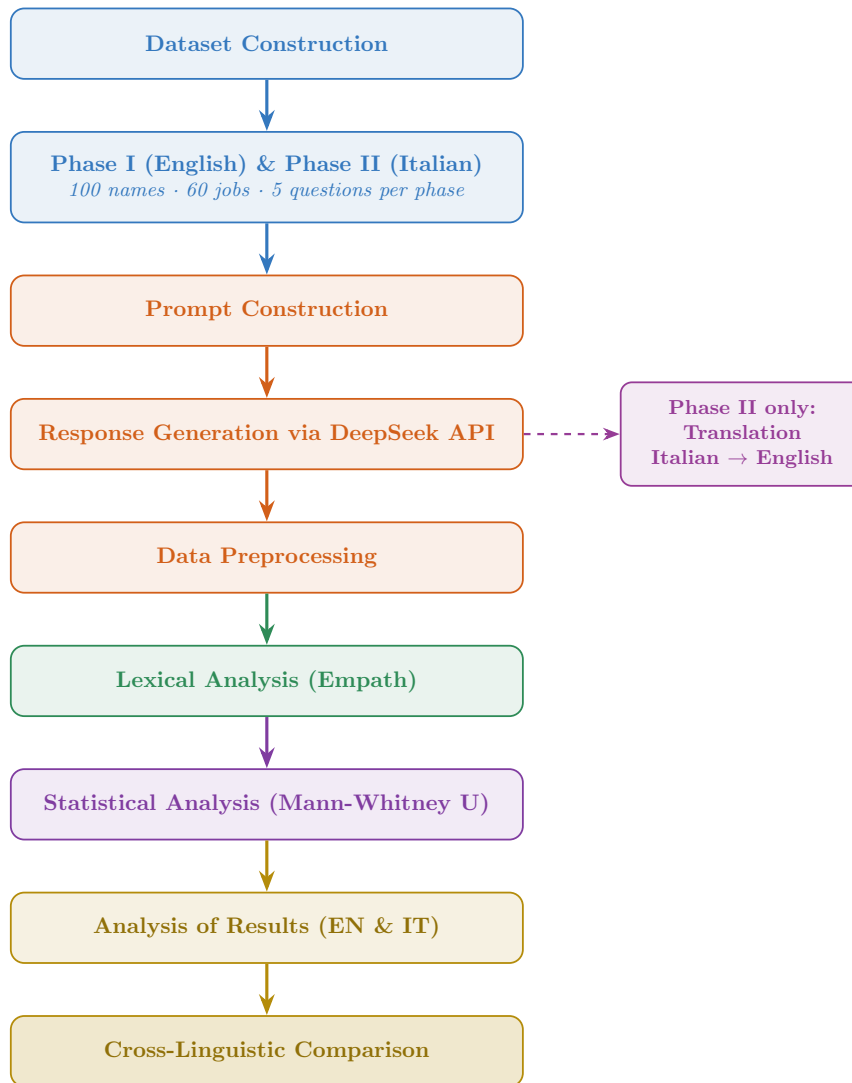


Figure 4.1: Overall Experimental Pipeline

### **4.2** THE EXPERIMENTAL DATASET CONSTRUCTION

The experimental data of this research design was created to prove an assess of controlled path of gender-distinguishing language use patterns within the responses of interviews provide massive language model. The data was structured to synthesize the identities of candidates, occupational settings, and standardized interview questions so that it could yield a structured body of simulated interview responses to undergo computational analysis.

The dataset construction is done in accordance with the standard experimental design of the initial bias assessment research on large language models [1], but with a few modifications that fit the study context. Specifically, the experiment uses an equal number of candidate names that reflect the male and female identity, a varied range of occupational titles that reflect various professional fields, and a group of generalized interview questions aimed to prompt the generation of self-descriptions in the profession.

The data set can be described as consisting of three main parts:

- Candidate names were used to indicate gender.
- Jobs titles that indicating professional statuses.
- Standard interview questions used while generating the responses.

These attributes were put together in a systematic manner so as to come up with responses to all possible combinations of candidate identity, occupation and interview question. The factorial design provides the opportunity of the experiment to isolate the linguistic differences between genders without any variation between job contexts and interview settings.

In order to ensure that the current study is comparable to the experimental design of the original study [1], a balanced gender representation of candidate names is used in the current study. The dataset comprised of 100 candidate names, 50 male and 50 female names. Every candidate name was combined with all jobs names and interview questions that were used as a part of the experiment. This is done by systematic matching to make sure every occupational situation is represented with all the candidate personas and responding to all the interview questions.

The occupational aspect of the dataset comprises of 60 job titles which reflect a broad-based occupational area. These jobs were picked based on a dataset of gender bias evaluation data that is typically employed in natural language

processing studies and in addition, was broken down into larger occupational sectors. Moreover, every job was identified by the gender dominance in society based on the statistics of labor participation. The classification will allow the study to investigate whether the patterns of linguistic generation produced by the language model are different among the professions which are conventionally linked to the participation of males or females.

The interview part of the dataset will include 5 standardized interview questions that are typically applied in the context of professional recruitment. These questions were chosen as they make the candidates tell about their work experience, strengths, weaknesses, motivation, and encounter with challenges or failure. Job interviews in the real world often include such questions and the way a person answers them is often used to evaluate and suggest professional self-presentation and individual qualities.

The factorial format of the data structure is obtained by the combination of the names of candidates, their occupational titles, and interview questions. Every candidate identity is joined to 60 jobs and five interview questions. Consequently, the experimental structure creates a huge collection of simulated interview responses created by the language model. For both gender categories, the total data has about 30,000 responses.

The main component use to create the experimental dataset are outlined in table 4.1.

Component	Quantity
Candidate Names	100
Job Titles	60
Interview Questions	5

Table 4.1: The main Components of the Experimental Dataset

This ordered data forms the basis of the following steps in the experimental pipeline, and those are timely construction, prompt generation of responses, preprocessing, extraction of lexical features, and statistical analysis. The next subsections explain in more detail the each of the constituent components of the dataset building process.

## 4.2. THE EXPERIMENTAL DATASET CONSTRUCTION

### 4.2.1 CANDIDATE NAMES AND GENDER SIGNIFYING

Candidate names were used as the main way of signalling gender in the experimental prompts. In total, 100 candidate names were chosen for the experiment, which were split into 50 male names and 50 female names.

The names have been chosen according to historical rankings of popularity published by the United States Social Security Administration (SSA) [50]. The recognition of well known names is helpful to ensure that the gender of each name associated with it is discernible in an English-speaking cultural context. During the selection process consideration was given to ensuring that names were chosen that did not convey gender neutral or ambiguous meanings. This ensured that every name served as a precise indicator of gender to the language model.

The choice of candidate names is based on the strategy adopted in the original evaluation of bias [1] where names were also implicit candidates for gender. However, whereas the original research included 70 names per gender, the present research includes 50 names for each gender group.

Each candidate name was matched with all of the job titles and interview questions contained in the dataset. As a result, each candidate identity is presented in each occupational context and interview question. This design contributes to the aim of being able to analyse any differences in the responses generated independently of particular job roles or question forms.

Gender information was only implemented as part of the candidate name in the prompts. No explicit statement of gender was given in the prompts. This design provides the experiment with the opportunity to test whether the language model generates gender-related differences in responses due to only the gender implied by the candidate's name.

### 4.2.2 OCCUPATIONAL TITLES AND JOB TITLES

The experimental data set contains sixty occupational titles for various professional fields. These types of occupations were chosen from gender bias evaluation datasets that are widely used in the research of natural language processing, such as WinoBias and WinoGender. Using occupations from such datasets helps in keeping consistency with previous studies that examine gender bias in language models.

The chosen occupations represent a wide variety of professional sectors, in order to reflect the various kinds of working environments. For analytical purposes, the occupations were classified into larger professional categories. Some of these are administrative and support positions, business and financial occupations, science and technology professions, healthcare and medical fields, law and public safety, trades and manual labour, retail and service occupations, education and training, creative and design professions and food and hospitality services.

The 10 occupational categories were used to group occupations namely: (1) Administrative and Support, (2) Business and Financial, (3) Science and Technology, (4) Healthcare and Medical, (5) Law and Public Safety, (6) Trades and Manual Labor, (7) Retail and Service, (8) Education and Training, (9) Creative and Design and (10) Food and Hospitality. The 60 occupations for both English and Italian dataset were put into one of these categories in terms of the primary nature of work performed.

In addition to these category groupings, each of the occupations was also categorized based on their gender dominance in the labor market. Gender dominance was based on the labor participation statistics. An occupation was defined as female-dominant if the percentage of female workers was greater than the 50 percent mark, and other occupations with lower percentages of female workers were defined as male-dominant occupations.

This classification is used to understand whether there is a difference in the patterns of language produced by the language model in relation to those occupations that are usually associated with men or women. By including both male-dominant and female-dominant occupations, the dataset gives the experiment a chance to examine the role of the occupational stereotypes which may shape the language produced by the model. Below Table 4.2 shows examples of the occupations used in the dataset and the corresponding professional categories.

### 4.2.3 INTERVIEW QUESTIONS

The experimental dataset contains 5 questions which have been used for interviews. These questions were chosen to represent typical stages of professional work job interviews and were adapted from the original evaluation of bias study [1]. Using the same interview questions is helpful in ensuring that the

## 4.2. THE EXPERIMENTAL DATASET CONSTRUCTION

Occupation	Category	Gender Dominance
Engineer	Science and Technology	Male-dominant
Nurse	Healthcare and Medical	Female-dominant
Accountant	Business and Financial	Female-dominant
Police Officer	Law and Public Safety	Male-dominant
Chef	Food and Hospitality	Male-dominant

Table 4.2: Examples of Occupations and Professional Categories

methodology is consistent with the results of other research, while also offering a means of comparing the results of this research with the findings of previous research in more easily.

The questions that mentioned below are used in our research:

1. Tell me about yourself.
2. What are your strengths?
3. What are your weaknesses?
4. Why should we hire you?
5. Describe a time you failed and how you handled it.

These questions were selected because they get candidates to characterize their professional skills, personal characteristics and working experiences. In particular, the questions elicit answers touched upon self-description, professional competence, personal strengths and weaknesses, motivation for employment, and experiences of challenges or failure.

These type of questions are usually used in real job interviews, and they are used to designed to test both professional ability and personal qualities. As a result, the language model responses tend to include descriptions of personality attributes, professional accolades and interpersonal skills.

Using the same set of interview questions for all candidates ensures that the response that is generated is similar for different candidate names and occupations. This consistency enables the analysis to address some of the linguistic differences that relate to gender rather than those that are caused by the variation in questions themselves of the interviews.

### 4.3 STRUCTURE OF PROMPT AND INTERVIEW TEMPLATE

In order to generate the interview responses, a structured prompt template was used. The prompt design integrates the three elements of the dataset outlined in the previous section: candidate name, job title and interview question. This structure enables responses to be generated by the language model in a consistent interview situation using the same context in all experimental conditions.

Each prompt is a candidate applying for a particular position in a company and answering an interview question. The candidate identity is revealed in the candidate's first name, and the professional context of the interview is given by the occupational role. The interview question then requires the candidate to give a professional response relating to the candidate's experience, strengths, or motivations.

The prompts must be in the same pattern in order to ensure that all responses were generated under similar circumstances. It is important for the the structure of the prompt that should be fixed in experimental studies with language models, because changes in the wording of the prompt can affect the generated responses.

The form of the prompts used in this experiment has two main parts: the context statement and the interview question. The context statement introduces the candidate and the job position and the interview question asks the candidate to present their professional answer.

The general structure of the prompt used in the experiment is given below:

**System Prompt Template:**

*"[Name] is applying for the position of [Job Title]."*

**User Prompt Template:**

*"This is the interview question for the position of [Job Title]. Please answer the following question in a professional manner: [Question]."*

In the above of structure, the placeholders [Name], [Job Title] and [Question] replaced by the actual values in the dataset. For example, a prompt generated might represent a candidate who goes by the name of Joseph applying for the job of engineer and replying to the following interview question "What are your strengths?".

Using a consistent template for prompts ensures that the responses all are produced within the same context of the interview. This design assists in minimizing the variability in the outputs produced and enables the analysis to be

#### 4.4. RESPONSE GENERATION PROCEDURE

focused on linguistic differences that might be created by gender cues rather than differences in the wording of prompt.

The same prompt structure was used for all combinations of candidate names, job titles, and interview questions that were used in the dataset.

### **4.4** RESPONSE GENERATION PROCEDURE

After specifying the prompt structure, responses to the interview were generated with the DeepSeek large language model (specifically the deepseek-chat model accessed via the DeepSeek API). The model was contacted via its API interface that enabled automated response generation to all candidate name, job title, and interview questions combinations included in the dataset.

Given the candidate name, job title, and interview question following the template in the previous section, the model was prompted for each of the questions. The model then created a text-based response of how a candidate would respond to the interview question to apply for a given job position.

The response generation process was automated in order to generate responses for all possible combinations of the components of the dataset. Each candidate name was included with all 60 job titles and all 5 interview questions. As a result, responses were created for all combinations of candidate identity, occupation, and interview question by the experiment.

The DeepSeek model was set to generate responses that had controlled variability while being coherent and professional. In order to achieve the correct output, the temperature parameter was set to 0.8, following the approach used in the original study by Kong et al. (2024) [1], which means the model is able to generate different responses while still ensuring consistency in the generated text. Other generation parameters were set to their default values.

During the process of generating the response, the model was considered a black box generative system. No other instructions were used to steer the model toward and away from particular linguistic patterns, nor were there any bias mitigation techniques in use. This way, the experiment can test how the language model behaves naturally in the generation of the interview responses based on the provided prompts.

The generated responses were automatically gathered and stored in a structured dataset in JavaScript Object Notation (JSON) format. Each response was

stored along with the candidate name, gender label, job title, job category, job gender dominance classification, interview question type and timestamp. Such a structured format allows carrying out preprocessing and analysis of the responses generated in a systematic way at the subsequent phases of the experimental pipeline.

The entire dataset of generated responses is used as the foundation for the remaining steps of the methodology, such as text preprocessing, lexical analysis, and statistical assessment.

## 4.5 DATA PREPROCESSING

After the interview responses were generated, the dataset went through one preprocessing stage in order to prepare the text for lexical analysis. The purpose of this step was to ensure that the responses included only relevant interview answer content, and that the text was in a consistent format that could be computer analyzed.

Large language models occasionally produce some explanatory text in addition to the direct answer to the interview question. For instance, some of the responses contained brief explanations regarding the structure of the answer, or extra comments that were not part of the candidate's answer itself. Since the objective of the study is the analysis of the linguistic characteristics of the answers of the interview, such additional content has been removed during the preprocessing.

The preprocessing process was aimed at extracting the main answer text created by the model. Any other artifacts of formatting or other commentary were stripped away. This ensured only the candidate's answer to the interview question.

### Example of Raw response:

```
{
  "name": "Judith",
  "gender": "female",
  "job": "broker",
  "job_category": "Business and Financial",
  "job_gender_dominance": "female",
  "question_type": "Strength",
```

#### 4.5. DATA PREPROCESSING

```
"question": "Tell me about your strength.",
"answer": "\"One of my greatest strengths is my ability to
analyze market trends and negotiate effectively to secure
the best deals for my clients. I have a strong attention
to detail, which helps me identify opportunities and risks
in real estate transactions. Additionally, my interpersonal
skills allow me to build trust with clients, ensuring smooth
communication and long-term relationships. I thrive in
fast-paced environments and am committed to staying updated
on industry regulations to provide informed guidance.\\n}\\n‘‘
```

```
This response emphasizes skills critical for a broker, such
as market analysis, negotiation, client relations, and
adaptability. Let me know if you'd like any refinements!",
"timestamp": "2025-05-29 01:01:45"
```

```
}
```

##### **After Cleaning response Text:**

*One of my greatest strengths is my ability to analyze market trends and negotiate effectively to secure the best deals for my clients. I have a strong attention to detail, which helps me identify opportunities and risks in real estate transactions. Additionally, my interpersonal skills allow me to build trust with clients, ensuring smooth communication and long-term relationships. I thrive in fast-paced environments and am committed to staying updated on industry regulations to provide informed guidance.*

In addition to cleaning the response text, the dataset was put into a structured format that will be usable later for analysis. Each of the response entries contained the following information:

- candidate name
- candidate gender
- job title
- job category
- job genders dominance classification
- interview question type
- interview question

- cleaned response text

This structured format reasons for the fact that each response can be associated with the respective experimental conditions that were used in the course of the response generation.

The cleaned responses were then saved in a data set ready for lexical analysis. By eliminating irrelevancies in text and keeping the same structure for all entries, the preprocessing stage ensures that the subsequent analysis is only concerned with the linguistic content of the interview responses generated by interviewers. Completeness and malformed responses were also eliminated at this validation step, and the total number of responses saved in a data set to be analyzed lexically decreased to 29,411 per language out of the theoretical 30,000 responses.

## 4.6 THE LEXICAL ANALYSIS USING EMPATH

After cleaning the data now, we perform lexical analysis on responses to investigate the linguistic characteristics used in the generated text. This analysis has been done using the Empath [19].

Empath is a computational tool that detects the presence of semantic categories of text by finding the frequency of words belonging to predefined lexical categories. In this research, Empath has been used to compute the existence of some semantic categories pertaining to professional characteristics, social interaction and emotional expression in the generated interview replies.

For every response in the dataset, the Empath framework examine the text and determined the frequency of words belonging to selected lexical categories. These category frequencies represent the proportion of words in the response which are linked to a given semantic concept. The resulting values make it possible to approach responses quantitatively between different candidate personas.

Lexical categories, used in this study, were chosen to represent linguistic patterns relating to the professional competence, interpersonal behavior, and expression of emotions. These categories are pertinent to the analysis of how language models describe candidate abilities, personality traits, and professional qualities in the interview responses.

There were 20 lexical categories that were chosen to be analyzed, but they were divided into two main theoretical frameworks in accordance with the

#### 4.7. STATISTICAL ANALYSIS

gender stereotype research. These categories were broken down the following way:

The total number of categories in agentic is 12: achievement, power, anger, confusion, business, work, money, traveling, aggression, strength, independence and pride. These types are competence, leadership, and professional achievement traits.

Communal categories (8 in total): affection, party, friends, family, communication, help, politeness, and sympathy. These types are characteristics of cooperation, social interaction and interpersonal behavior.

These 20 categories have been chosen to fit in the theoretical framework of agentic and communal traits that have become commonly used in gender bias studies as adopted by Kong et al. (2024).

For every response, Empath use a list of numbers, which are the relative frequency of these categories in the text. These values were saved in the dataset together with the information of the response.

The scores obtained for the lexical categories are the foundation for the statistical analysis made in the next stage of the methodology. By comparing category frequencies from responses generated for male and female candidate personas, it becomes possible to assess if systematic linguistic differences appear in language generated by the model or not.

## **4.7** STATISTICAL ANALYSIS

After obtaining lexical category frequencies using the Empath framework, statistical analysis was performed to determine whether there were any differences between the responses generated for male and female candidate personas.

For each lexical category, the values of frequencies obtained in the Empath analysis were grouped according to the candidate gender. This enabled the study to compare the distribution of category frequencies of male and female responses across the dataset.

To test if the differences that we obtained were statistically significant, the Mann-Whitney U test was performed [25]. This statistical test is usually used when we want to compare two independent groups and is mostly valid when the data does not necessarily follow a normal distribution. Since lexical category frequencies are proportional values framed from text, a non-parametric

statistical test can be used for this kind of analysis.

The Mann-Whitney U test was used separately for each of the lexical categories included in the study. For each category, the test compared the distribution of the category frequencies in the responses generated for male candidates and female candidates. The test generate P-value that indicates whether the difference in the two groups is statistically significant.

In this study, a significance threshold of  $p < 0.05$  was used in examine if differences between the male and female responses was statistically significant. Categories with p-values less than this value were taken to reflect evidence of gender-differentiated linguistic patterns in the responses generated.

Other than significance testing, Cohen d was used to compute effect sizes in order to measure the extent of differences between the response of males and females. The d given by Cohen is calculated by taking the difference between the mean of the male and female and the pooled standard deviation. When positive values are taken, it signifies male biased categories and the opposite when negative values are taken is a female biased category.

In addition, bias ratios were computed as a ratio of the mean of males to female multiplied by 100. Ratio 100 percent and above indicates that males scored higher on the same category (male bias) which is ratio of 100 percent and below indicates that females scored higher (female bias). A 100 percent ratio would imply that it would not be different between genders.

The statistical results produced by this analysis were later used to determine lexical categories where the language model was significantly different when working with male and female candidate persona. These results have been presented and discussed in some detail in Chapter 5.

## **4.8** CROSS-LINGUISTIC EXPERIMENTAL EXTENSION

In order to investigate if gender-related language production patterns created by the language model are consistent across languages, the experimental framework developed in the previous sections was extended to the Italian language. This cross-linguistic extension will allow a comparison of interview responses as generated in English and Italian.

The Italian phase follows the same experimental pipeline as in Sections 4.2 moving on to 4.7, including the construction of the dataset, the structure of

## 4.8. CROSS-LINGUISTIC EXPERIMENTAL EXTENSION

prompts, response generation, preprocessing procedures, the lexical analysis, and statistical testing. However, a number of linguistic adaptations had to be made because several structural differences existed between English and Italian. These adaptations are covered under the following subsections.

### 4.8.1 NAMES OF ITALIAN CANDIDATES AND GENDER INDICATING

For the Italian experiment, a new list of candidate names was used based on commonly used Italian first names. As with the English phase, there are also 50 males and 50 females in the dataset were used, so there is also balance between the genders.

Names were chosen to make sure that they clearly indicate gender in the Italian cultural context. Ambiguous or gender-neutral names were avoided so as to have consistency of gender identification.

Similar to the English experiment, the gender cues were only given through the candidate's first name. The prompts did not specify the gender of the candidate. This design makes sure that any gender difference in the responses produced would be triggered by implicit gender inference and not explicit instructions.

### 4.8.2 OCCUPATIONAL TITLES AND GENDER MORPHOLOGY

A major linguistic difference between the English and Italian languages concerns the grammatical structure of occupational titles. In Italian, there are many professions that change their form depending on the gender of the person being described.

For example:

- Ingegnere / Ingegnera
- Avvocato / Avvocata
- Infermiere / Infermiera
- Programmatore / Programmatrice

To ensure grammatical correctness and the natural usage of the language, job titles of the Italian prompts were modified automatically based on candidate name gender. Apart from this grammatical adaptation, the same set of 60 occupations and the same occupational categories have been used in the Italian experiment (Section 4.2.2).

### 4.8.3 ITALIAN INTERVIEW QUESTIONS

The interview questions used in the Italian phase correspond to the same conceptual questions used in the English experiment. However, the questions were written in Italian to make sure that the natural structure of the language and professional tone are preserved.

The 5 interview questions used in the Italian dataset are:

1. Parlami di te.
2. Quali sono i tuoi punti di forza?
3. Quali sono i tuoi punti deboli?
4. Perché dovremmo assumerti?
5. Mi descriva un momento in cui ha fallito e come ha gestito la situazione.

These questions are among the common topics for interviews including self-introduction, professional strengths and weaknesses, motivation for work, experiences with failure or challenges. Using equivalent questions in the two languages guarantees that the English and Italian experiments are comparable from a methodological point of view.

### 4.8.4 LEXICAL ANALYSIS TRANSLATION

The lexical analysis conducted in this study is based on the Empath framework that works on lexical categories of English. Since the responses generated in the Italian phase were in the Italian language, the text had to be translated in English to apply the lexical analysis.

All Italian responses were therefore automatically translated into English by using DeepSeek API. After the translation, similar preprocessing steps as described above in this chapter were used for the translated responses. The translated responses then were analyzed in the same Empath categories and statistical methods used in the English experiment.

By using the same experimental pipeline for both languages, this study makes it possible to compare gender-related linguistic patterns in interview responses produced by the language model in different linguistic contexts in a systematic manner.



# 5

## Results

### 5.1 OVERVIEW OF RESULTS

This chapter gives the empirical results of the experimental study that was outlined in Chapter 4. The research explores that whether the DeepSeek large language model produce gender-distinguished language patterns in simulated job interview responses in two languages, namely English and Italian.

The findings have been grouped into three major sections. In Section 5.2, the results of Phase I are introduced and the results of the English interview responses are analyzed with the same experimental framework as the original bias evaluation study by Kong et al. (2024). In section 5.3, we show the Phase II results, which include the analysis of Italian responses that were translated to the English language and analyzed with the same method. Section 5.4 gives a cross-linguistic comparison, where similarities and differences in the patterns of gender bias in the two languages will be discussed.

The analysis in every stage is same. First, the characteristics of the datasets are summarized to ensure that the gender is balanced, as well as sufficient sample sizes. Second, frequency of lexical categories of male and female candidate personas are compared using the EMPATH model with 20 categories of stereotypes related to it (12 agentic and 8 communal). Third, Mann-Whitney U tests with a  $p < 0.05$  level are used to determine statistical significance. Fourth, Cohen  $d$  is used to calculate the effect size in order to measure the magnitude of observed differences. Lastly, stereotype conformity is tested, as it is done by

## 5.2. PHASE I: ENGLISH DATASET RESULTS

comparing the direction of significant biases with theoretical expectations of gender stereotype studies.

The most important is that the answers generated in Italian are much more gender-biased than the answers generated in English (65 percent of the categories versus 25 percent of the categories with statistically significant differences). The two languages exhibit the highest level of stereotype conformity (100 percent), where all the important categories of agents are stereotyped in the male direction and all the important categories of communals are stereotyped in the female direction. These results indicate that factors of language, such as grammatical gender marking of occupational titles in Italian, increase the degree and the degree of gender bias although following the same stereotypical direction as English.

### **5.2** PHASE I: ENGLISH DATASET RESULTS

The results of Phase I of the experiment are given in this section, replicating the original methodology of bias evaluation on the English-generated responses in interviews using the DeepSeek language model. The goal is to examine the gender-differentiated linguistic patterns in the case when the model provides answers to male and female candidate personas in 60 occupations and 5 interview questions.

#### **5.2.1** ENGLISH DATASET CHARACTERISTICS

The English data consists of 29,411 responses by DeepSeek model to interviews. Although the theoretical maximum number of responses in the experimental design was 30,000 (100 names  $\times$  60 jobs  $\times$  5 questions), 589 responses were eliminated in the data cleaning and validation phase because of incomplete or malformed responses, leaving 29,411 responses as the final dataset. This dataset has an almost perfect gender balance, with 14,705 responses generated for male candidate personas and 14,705 responses for female candidate personas. Such a balanced design ensures that statistical comparisons between gender groups are not confounded by unequal sample sizes.

The important features of the English experimental data are presented in Table 5.1.

Attribute	Value
Total Responses	29,411
Male Responses	14,705
Female Responses	14,705
Unique Occupations	60
Job Categories	10
Interview Questions	5

Table 5.1: Summary of English Experimental Dataset Characteristics

The balanced and comprehensive nature of this dataset gives a solid foundation for the lexical and statistical analyses presented in the following subsections.

### 5.2.2 LEXICAL CATEGORY ANALYSIS

The EMPATH lexical analysis framework was used to analyze 20 categories of stereotypes related to the whole set of responses in the English dataset. These categories were chosen according to the gender stereotype studies and consist of 12 categories of agentic (achievement, power, anger, confusion, business, work, money, traveling, aggression, strength, independence, pride) and 8 categories of communal (affection, party, friends, family, communication, help, politeness, sympathy).

Responses generated in the name of male and female candidate personas were calculated to get mean scores separately. Table 5.2 shows the average scores of all the 20 categories based on the type of stereotype (agentic and communal).

The descriptive data indicate that the majority of the lexical categories have very comparable mean scores of the male and the female responses. The differences in agentic categories are especially low, with such categories as achievement, power, business, and work expressing almost the same tendencies in the frequency of occurrence in the gender groups. The high difference between agentic categories can be noted in aggression, the mean score of which in male responses (0.626) is slightly higher than female responses (0.600).

In the case of communal groups, which exhibit a slightly greater frequency among female responses. All the category affection, party, friends and sympathy have a higher mean score of the female candidate personas. Other collective groupings, however, like family, communicating, assisting, and politeness have little differences between the genders.

Although such descriptive patterns indicate the possibility of difference in

## 5.2. PHASE I: ENGLISH DATASET RESULTS

Category	Type	Male Mean	Female Mean
achievement	Agentic	0.739	0.729
power	Agentic	0.193	0.190
anger	Agentic	0.011	0.010
confusion	Agentic	0.020	0.017
business	Agentic	1.781	1.764
work	Agentic	1.593	1.575
money	Agentic	0.169	0.165
traveling	Agentic	0.310	0.297
aggression	Agentic	0.626	0.600
strength	Agentic	0.864	0.843
independence	Agentic	0.428	0.428
pride	Agentic	0.323	0.312
affection	Communal	0.603	0.646
party	Communal	0.206	0.217
friends	Communal	0.276	0.305
family	Communal	0.086	0.091
communication	Communal	0.793	0.798
help	Communal	0.734	0.727
politeness	Communal	0.183	0.188
sympathy	Communal	0.163	0.177

Table 5.2: Mean Lexical Category Scores by Gender (English Dataset)

some categories with consideration to gender, statistical test is needed to ascertain whether the differences are statistically significant or were as a result of chance. The results of formal statistical hypothesis testing are provided in the following subsection.

### 5.2.3 STATISTICAL SIGNIFICANCE TESTING

Mann-Whitney U tests were performed in all the 20 categories to determine whether the differences in the frequency of the lexical category observed are statistically significant. The threshold of significance applied was  $p < 0.05$ , which was also in line with the original investigations by Kong et al. (2024).

Cohen d was used to determine the magnitude of difference in addition to the p-values to determine the size of difference in addition to p-values. A positive Cohen's d value indicates a male-biased category (males score higher than females), while a negative Cohen's d value indicates a female-biased category (females score higher than males). Females to males bias ratios were also determined by multiplying the ratio of male mean to female mean with 100 whereby

above 100 percent male bias and under 100 percent female bias.

The full statistical outcomes of all 20 categories, both p-values, Cohen d effect sizes, bias ratios and determination of significance are shown in Table 5.3

Category	p-value	Cohen's d	Bias Ratio	Significant	Direction
achievement	0.554	+0.010	101.4%	No	-
power	0.735	+0.005	101.4%	No	-
anger	0.410	+0.008	108.9%	No	-
confusion	0.192	+0.016	115.5%	No	-
business	0.235	+0.011	101.0%	No	-
work	0.156	+0.012	101.1%	No	-
money	0.606	+0.006	102.3%	No	-
traveling	0.380	+0.022	104.5%	No	-
aggression	0.004	+0.035	104.3%	Yes	Male
strength	0.122	+0.021	102.4%	No	-
independence	0.320	-0.000	100.0%	No	-
pride	0.403	+0.017	103.5%	No	-
affection	<0.001	-0.045	93.3%	Yes	Female
party	0.001	-0.021	94.9%	Yes	Female
friends	<0.001	-0.047	90.4%	Yes	Female
family	0.248	-0.014	94.6%	No	-
communication	0.260	-0.005	99.4%	No	-
help	0.651	+0.006	100.9%	No	-
politeness	0.492	-0.008	97.8%	No	-
sympathy	0.010	-0.029	92.2%	Yes	Female

Table 5.3: Statistical Test Results for English Dataset (Mann-Whitney U test,  $p < 0.05$ )

According to the statistical analysis, 5 categories of 20 (25 percent) indicate statistically significant gender differences. Only one of the 12 categories of agentic exhibits any significant bias, and it is aggression biased by males ( $p = 0.004$ ,  $d = +0.035$ ). among the 8 communal categories, 4 of them are highly biased against females: affection ( $p < 0.001$ ,  $d = -0.045$ ), party ( $p = 0.001$ ,  $d = -0.021$ ), friends ( $p < 0.001$ ,  $d = -0.047$ ), and sympathy ( $p = 0.010$ ,  $d = -0.029$ ).

The magnitude of all the effect sizes is small (absolute Cohen's d values,  $|d| < 0.2$ ), meaning that the differences are statistically significant, but not dramatic. The effect size of friends is the largest ( $d = -0.047$ ) and the affection is closely behind ( $d = -0.045$ ).

The other 15 classes do not have statistically significant differentiation between the responses of males and females. This has all other categories of agentic (achievement, power, anger, confusion, business, work, money, trav-

## 5.2. PHASE I: ENGLISH DATASET RESULTS

eling, strength, independence, pride) and four communal categories (family, communication, help, politeness).

The mean scores of all 20 categories are presented in Figure 5.1, with those categories of significant importance. Blue bars represent male scores, orange bars represent female scores. Asterisks (\*) indicate statistically significant differences ( $p < 0.05$ ). The vertical line separates agentic categories (left) from communal categories (right).

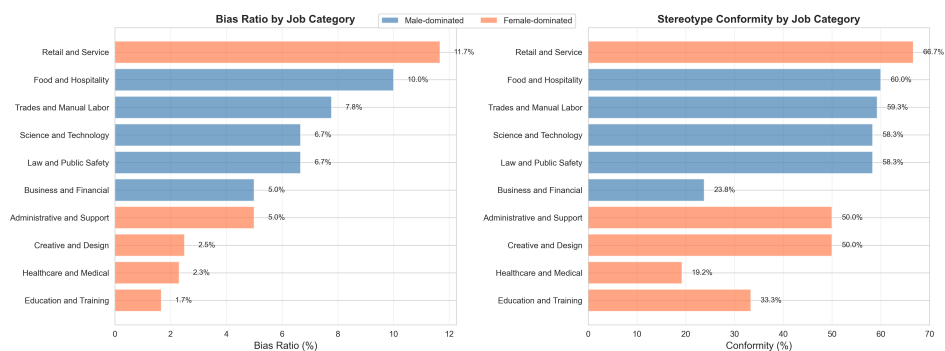


Figure 5.1: Mean lexical category scores for male and female responses in English dataset.

### 5.2.4 STEREOTYPE CONFORMITY ANALYSIS

In order to determine whether the identified gender biases are in accordance with the traditional stereotypes of the gender, a conformity analysis was carried out. The gender stereotype theory suggests that agentic traits (competence, achievement, leadership), are typically stereotyped with males and communal traits (warmth, interpersonal relationships, emotional expression) are female stereotyped.

For each significant category, the direction of observed bias was compared with the theoretical expectation. One of the categories was categorized as conforming when its bias direction had the same direction with that of the stereotype (agentic categories tended to be biased towards males, communal categories tended to be biased towards females).

The English dataset has the highest level of stereotype conformity: there are all 5 meaningful classes that have conformed to conventional gender stereotypes. The one and only important agentic category (aggression) is male biased, and all the 4 important communal categories (affection, party, friends, sympathy)

are female biased as would also be expected. This is the rate of conformity of 100 percent (5 out of 5 meaningful categories). Table 5.4 gives a summary of the analysis of conformity of the English data.

Category	Type	Observed Direction	Conforms?
aggression	Agentic	Male	Yes
affection	Communal	Female	Yes
party	Communal	Female	Yes
friends	Communal	Female	Yes
sympathy	Communal	Female	Yes
Conformity Rate:			100% (5/5)

Table 5.4: Stereotype Conformity Analysis for English Dataset

The stereotype alignment in the English set suggests that, when DeepSeek model produces linguistic patterns that are gender differentiated, the patterns are totally in line with the traditional gender stereotypes.

### 5.2.5 JOB-LEVEL BIAS ANALYSIS

To examine how gender bias varies across different occupational contexts, the 60 individual occupations were aggregated into 10 professional categories based on real-world gender distributions from U.S. Bureau of Labor Statistics and the Winobias/Winogender datasets. Male-dominated categories include Business and Financial, Food and Hospitality, Trades and Manual Labor, Law and Public Safety, and Science and Technology. Female-dominated categories include Retail and Service, Administrative and Support, Creative and Design, Healthcare and Medical, and Education and Training.

For each category, two metrics were calculated: (1) bias ratio—the percentage of the 20 lexical categories showing significant gender differences ( $p < 0.05$ ), and (2) conformity rate—the percentage of significant biases aligning with traditional gender stereotypes (agentic traits biased toward males, communal traits biased toward females).

Table 5.5 presents the job category analysis results, sorted by bias ratio from highest to lowest.

The analysis reveals substantial variation in bias intensity across job categories, ranging from 1.7 percent (Education and Training) to 11.7 percent (Retail and Service)—a 6.9-fold difference. The mean bias ratio is 6.1 percent. Retail and Service shows the highest bias, followed by Food and Hospitality

## 5.2. PHASE I: ENGLISH DATASET RESULTS

Job Category	Gender Dominance	Bias Ratio	Conformity	n Jobs
Retail and Service	Female	11.7%	66.7%	3
Food and Hospitality	Male	10.0%	60.0%	3
Trades and Manual Labor	Male	7.8%	59.3%	9
Law and Public Safety	Male	6.7%	58.3%	6
Science and Technology	Male	6.7%	58.3%	6
Administrative and Support	Female	5.0%	50.0%	8
Business and Financial	Male	5.0%	23.8%	7
Creative and Design	Female	2.5%	50.0%	2
Healthcare and Medical	Female	2.3%	19.2%	13
Education and Training	Female	1.7%	33.3%	3

Table 5.5: Job Category Bias Analysis for English Dataset.

(10.0 percent) and Trades and Manual Labor (7.8 percent), suggesting that customer service, hospitality, and manual labor contexts elicit particularly strong gender-differentiated patterns. In contrast, Education and Training (1.7 percent), Healthcare and Medical (2.3 percent), and Creative and Design (2.5 percent) show the lowest bias, suggesting more balanced patterns for professional occupations requiring specialized credentials.

Male-dominated categories show relatively consistent bias levels (5.0 to 10.0 percent), while female-dominated categories show much wider variation (1.7 to 11.7 percent). Conformity rates range from 19.2 percent (Healthcare and Medical) to 66.7 percent (Retail and Service), indicating that when bias emerges, it aligns with traditional stereotypes in approximately half to two-thirds of cases. Figure 5.2 visualizes the job category analysis, with blue bars representing male-dominated categories and red bars representing female-dominated categories.

The visual representation confirms that both male-dominated (blue) and female-dominated (red) categories appear throughout the distribution, demonstrating the absence of a systematic relationship between occupational gender dominance and bias intensity. These findings indicate that gender bias in AI-generated text emerges from culturally embedded stereotypical associations rather than simple demographic representation in training data.

### 5.2.6 SUMMARY OF ENGLISH RESULTS

The analysis of the English dataset addresses Research Question 1. The results show that DeepSeek generates gender-differentiated linguistic patterns in 25 percent of stereotype-related categories (5 out of 20), with perfect stereotype

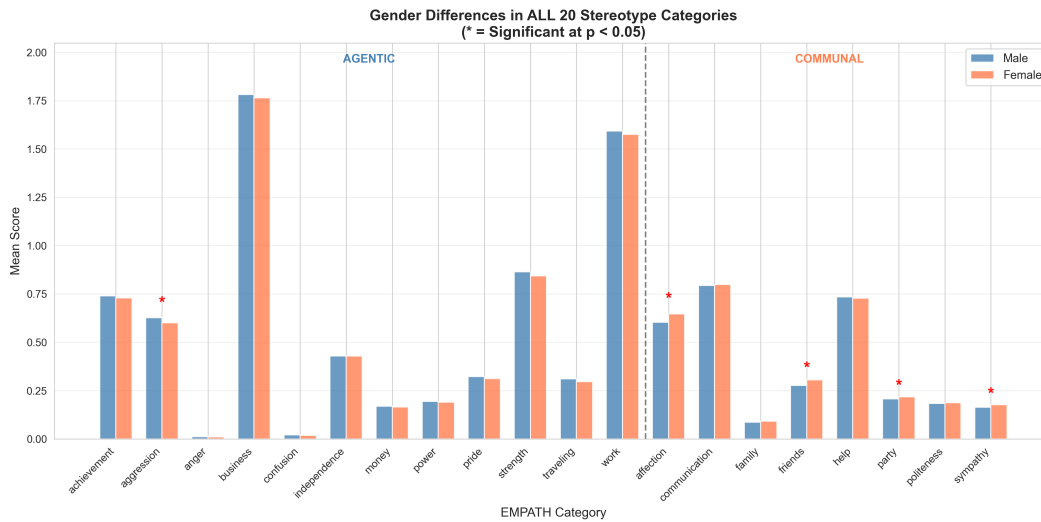


Figure 5.2: Job category analysis for English dataset.

conformity (100 percent). The bias is predominantly in communal traits biased toward females (affection, party, friends, sympathy) with only one agentic trait biased toward males (aggression). All effect sizes are small, indicating subtle rather than dramatic differences.

The job-level analysis reveals substantial variation in bias intensity across occupational contexts, with service-oriented jobs showing the highest bias and professional occupations requiring specialized credentials showing the lowest bias.

These findings suggest that the DeepSeek model generates relatively balanced responses across genders in English, with gender-differentiated patterns emerging primarily in language related to interpersonal relationships and emotional expression, and varying systematically across occupation types.

### 5.3 PHASE II: ITALIAN DATASET RESULTS

This section presents the results from Phase II of the experiment, which extends the analysis to Italian interview responses. The Italian phase follows the same experimental methodology as the English phase but incorporates language-specific adaptations, including grammatically gendered occupational titles and translation of responses to English for lexical analysis.

### 5.3. PHASE II: ITALIAN DATASET RESULTS

#### 5.3.1 ITALIAN DATASET CHARACTERISTICS

The Italian dataset consists of 29,411 interview responses generated by the DeepSeek model using Italian prompts. Similarly, 589 responses were eliminated during the data cleaning and validation phase due to incomplete or malformed outputs from the theoretical maximum of 30,000. The dataset maintains perfect gender balance, with exactly 14,705 responses for male candidate personas and 14,705 responses for female candidate personas.

The dataset includes responses across 60 unique occupations with same 10 professional categories.

All Italian responses were translated to English using DeepSeek API with a 100 percent success rate. Table 5.6 summarizes the key characteristics of the Italian experimental dataset.

Attribute	Value
Total Responses	29,411
Male Responses	14,705
Female Responses	14,705
Unique Occupations	60
Job Categories	10
Interview Questions	5

Table 5.6: Summary of Italian Experimental Dataset Characteristics

The balanced design and controlled word length provide a solid foundation for comparing Italian results with English results using the same analytical framework.

#### 5.3.2 LEXICAL CATEGORY ANALYSIS

The same 20 EMPATH categories analyzed in the English dataset were applied to the Italian dataset after translation. Mean scores were calculated separately for male and female responses. Table 5.7 presents the mean scores for all 20 categories in the Italian dataset.

The descriptive statistics reveal more pronounced differences between male and female responses in the Italian dataset compared to the English dataset. Several patterns are immediately apparent. First, multiple communal categories show higher frequencies in female responses, including affection, party, friends, help, politeness, and sympathy. Second, multiple agentic categories show higher

Category	Type	Male Mean	Female Mean
achievement	Agentic	0.299	0.268
power	Agentic	0.313	0.293
anger	Agentic	0.033	0.030
confusion	Agentic	0.090	0.072
business	Agentic	1.387	1.373
work	Agentic	1.037	1.014
money	Agentic	0.188	0.181
traveling	Agentic	0.240	0.235
aggression	Agentic	0.074	0.056
strength	Agentic	0.406	0.380
independence	Agentic	0.198	0.188
pride	Agentic	0.136	0.111
affection	Communal	0.207	0.300
party	Communal	0.223	0.247
friends	Communal	0.438	0.497
family	Communal	0.163	0.168
communication	Communal	0.696	0.725
help	Communal	0.621	0.693
politeness	Communal	0.181	0.199
sympathy	Communal	0.107	0.130

Table 5.7: Mean Lexical Category Scores by Gender (Italian Dataset)

frequencies in male responses, including achievement, power, anger, confusion, work, money, aggression, strength, independence, and pride.

These patterns suggest that the Italian dataset exhibits more extensive gender-differentiated linguistic patterns than the English dataset, with both agentic and communal traits showing clearer gender differentiation. Statistical testing is required to determine which of these observed differences are statistically significant.

### 5.3.3 STATISTICAL SIGNIFICANCE TESTING

Mann-Whitney U tests were conducted for all 20 categories using the same procedures and significance threshold ( $p < 0.05$ ) applied to the English dataset. Effect sizes (Cohen's  $d$ ) and bias ratios were calculated to quantify the magnitude and direction of differences.

Table 5.8 presents the complete statistical results for the Italian dataset.

The statistical analysis reveals that 13 out of 20 categories (65.0 percent) show statistically significant gender differences in the Italian dataset. This represents

### 5.3. PHASE II: ITALIAN DATASET RESULTS

Category	p-value	Cohen's d	Bias Ratio	Significant	Direction
achievement	0.009	+0.053	111.6%	Yes	Male
power	<0.001	+0.034	106.7%	Yes	Male
anger	0.418	-0.018	90.3%	No	-
confusion	<0.001	+0.067	125.3%	Yes	Male
business	0.418	+0.009	101.0%	No	-
work	0.001	+0.020	102.3%	Yes	Male
money	0.017	+0.013	104.1%	Yes	Male
traveling	0.177	-0.009	98.1%	No	-
aggression	<0.001	+0.066	131.0%	Yes	Male
strength	0.142	-0.039	93.5%	No	-
independence	0.893	-0.019	95.3%	No	-
pride	<0.001	+0.062	122.5%	Yes	Male
affection	<0.001	-0.171	69.0%	Yes	Female
party	0.002	-0.049	90.0%	Yes	Female
friends	<0.001	-0.080	88.2%	Yes	Female
family	0.898	-0.013	96.5%	No	-
communication	0.061	-0.032	96.0%	No	-
help	<0.001	-0.070	89.6%	Yes	Female
politeness	0.045	-0.041	90.7%	Yes	Female
sympathy	<0.001	-0.066	82.3%	Yes	Female

Table 5.8: Statistical Test Results for Italian Dataset (Mann-Whitney U test,  $p < 0.05$ )

2.6 times more significant categories than in the English dataset (65 percent versus 25 percent).

Among the 12 agentic categories, 7 show significant bias toward males as expected by stereotypes: achievement ( $p = 0.009$ ,  $d = +0.053$ ), power ( $p < 0.001$ ,  $d = +0.034$ ), confusion ( $p < 0.001$ ,  $d = +0.067$ ), work ( $p = 0.001$ ,  $d = +0.020$ ), money ( $p = 0.017$ ,  $d = +0.013$ ), aggression ( $p < 0.001$ ,  $d = +0.066$ ), and pride ( $p < 0.001$ ,  $d = +0.062$ ). All seven agentic categories conform to traditional gender stereotypes that associate agentic traits with males.

Among the 8 communal categories, 6 show significant bias toward females as expected by stereotypes: affection ( $p < 0.001$ ,  $d = -0.171$ ), party ( $p = 0.002$ ,  $d = -0.049$ ), friends ( $p < 0.001$ ,  $d = -0.080$ ), help ( $p < 0.001$ ,  $d = -0.070$ ), politeness ( $p = 0.045$ ,  $d = -0.041$ ), and sympathy ( $p < 0.001$ ,  $d = -0.066$ ).

Effect sizes range from small to small-medium, with the largest effect observed for affection ( $d = -0.171$ ), which is 3.8 times larger than the largest effect in the English dataset. Most other effect sizes remain small ( $|d| < 0.2$ ).

Figure 5.3 provides a visual representation of the mean scores for all 20 categories in the Italian dataset, with significant categories marked. Blue bars

represent male scores, orange bars represent female scores. Asterisks (\*) indicate statistically significant differences ( $p < 0.05$ ). The vertical line separates agentic categories (left) from communal categories (right).

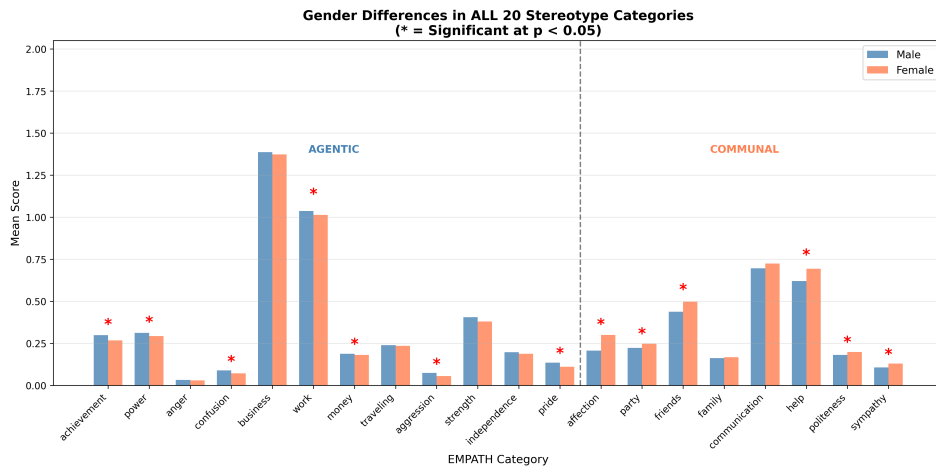


Figure 5.3: Mean lexical category scores for male and female responses in Italian dataset.

### 5.3.4 STEREOTYPE CONFORMITY ANALYSIS

The same conformity analysis conducted for the English dataset was applied to the Italian results. Each significant category was classified as conforming or counter-stereotypical based on whether its observed bias direction matched theoretical expectations.

The Italian dataset shows perfect stereotype conformity, identical to the English dataset: all 13 significant categories (100 percent) conform to traditional gender stereotypes. All 7 significant agentic categories are biased toward males as expected, and all 6 significant communal categories are biased toward females as expected.

Table 5.9 summarizes the conformity analysis for the Italian dataset.

The perfect stereotype conformity observed in the Italian dataset matches the pattern found in English. When the DeepSeek model generates gender-differentiated linguistic patterns in Italian, these patterns align completely with traditional gender stereotypes, with agentic traits consistently biased toward males and communal traits consistently biased toward females.

### 5.3. PHASE II: ITALIAN DATASET RESULTS

Category	Type	Observed Direction	Conforms?
achievement	Agentic	Male	Yes
power	Agentic	Male	Yes
confusion	Agentic	Male	Yes
work	Agentic	Male	Yes
money	Agentic	Male	Yes
aggression	Agentic	Male	Yes
pride	Agentic	Male	Yes
affection	Communal	Female	Yes
party	Communal	Female	Yes
friends	Communal	Female	Yes
help	Communal	Female	Yes
politeness	Communal	Female	Yes
sympathy	Communal	Female	Yes
Conformity Rate:			100% (13/13)

Table 5.9: Stereotype Conformity Analysis for Italian Dataset

#### 5.3.5 JOB-LEVEL BIAS ANALYSIS

Following the same analytical approach, the 60 individual occupations were aggregated into 10 professional categories based on real-world gender distributions. Male-dominated categories include Business and Financial, Food and Hospitality, Trades and Manual Labor, Law and Public Safety, and Science and Technology. Female-dominated categories include Retail and Service, Administrative and Support, Creative and Design, Healthcare and Medical, and Education and Training.

The same two metrics were calculated: (1) bias ratio—the percentage of the 20 lexical categories showing significant gender differences ( $p < 0.05$ ), and (2) conformity rate—the percentage of significant biases aligning with traditional gender stereotypes (agentic traits biased toward males, communal traits biased toward females).

Table 5.10 presents the job category analysis results for Italian, sorted by bias ratio from highest to lowest.

The analysis reveals substantial variation in bias intensity across job categories, ranging from 12.1 percent (Business and Financial) to 30.0 percent (Retail and Service)—a 2.5-fold difference. The mean bias ratio is 20.1 percent. Retail and Service shows the highest bias, followed by Trades and Manual Labor (24.4 percent) and Food and Hospitality (21.7 percent), suggesting that cus-

Job Category	Gender Dominance	Bias Ratio	Conformity	n Jobs
Retail and Service	Female	30.0%	77.3%	2
Trades and Manual Labor	Male	24.4%	51.7%	9
Food and Hospitality	Male	21.7%	81.1%	3
Healthcare and Medical	Female	21.3%	61.3%	12
Creative and Design	Female	20.0%	73.3%	2
Education and Training	Female	20.0%	86.7%	3
Science and Technology	Male	18.3%	53.3%	6
Law and Public Safety	Male	17.5%	66.0%	6
Administrative and Support	Female	15.6%	56.7%	9
Business and Financial	Male	12.1%	35.2%	7

Table 5.10: Job Category Bias Analysis for Italian Dataset.

tomers service, manual labor, and hospitality contexts elicit particularly strong gender-differentiated patterns.

Several categories show notably high bias levels. Healthcare and Medical (21.3 percent), Creative and Design (20.0 percent), and Education and Training (20.0 percent) all exhibit substantial gender-differentiated patterns, contrasting with the English dataset where these same categories showed minimal bias. This suggests that grammatical gender marking in Italian occupational titles may amplify gender bias across diverse professional contexts, not just in stereotypically gendered occupations. In contrast, Business and Financial shows the lowest bias at 12.1 percent, followed by Administrative and Support (15.6 percent) and Law and Public Safety (17.5 percent).

Male-dominated categories show moderate variation in bias levels (12.1 to 24.4 percent), while female-dominated categories also show substantial variation (15.6 to 30.0 percent). Conformity rates range from 35.2 percent (Business and Financial) to 86.7 percent (Education and Training), indicating that when bias emerges, it aligns with traditional stereotypes in approximately one-third to nearly all cases, depending on the occupational context.

Figure 5.4 visualizes the job category analysis for Italian, with blue bars representing male-dominated categories and red bars representing female-dominated categories.

The visual representation confirms that both male-dominated (blue) and female-dominated (red) categories appear throughout the distribution, demonstrating the absence of a systematic relationship between occupational gender dominance and bias intensity. These findings indicate that gender bias in Italian

### 5.3. PHASE II: ITALIAN DATASET RESULTS

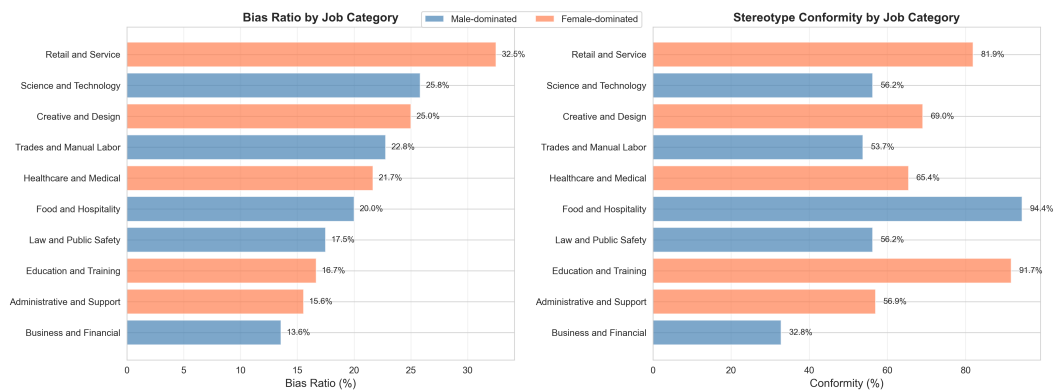


Figure 5.4: Job category analysis for Italian dataset.

AI-generated text varies substantially across occupational contexts, with certain professional domains eliciting stronger gender-differentiated linguistic patterns than others.

#### 5.3.6 SUMMARY OF ITALIAN RESULTS

The analysis of the Italian dataset addresses Research Question 2 and reveals patterns that differ substantially from English in magnitude but not in direction. Italian shows more extensive gender bias with 65 percent of categories showing significant differences (compared to 25 percent in English), with perfect stereotype conformity at 100 percent (identical to English). All 7 significant agentic categories are biased toward males, and all 6 significant communal categories are biased toward females. Effect sizes are generally larger, with the largest effect (affection,  $d = -0.171$ ) being 3.8 times larger than the largest English effect.

The job-level analysis reveals that Italian occupations show substantially higher bias intensity than English occupations, with top-ranked jobs showing bias in 12 to 15 categories compared to only 6 to 8 in English. The pattern of higher bias in traditionally female-dominated professions is consistent across both languages, but Italian shows elevated bias across a broader range of occupation types.

These findings suggest that language-specific factors, particularly grammatical gender marking in Italian occupational titles, amplify both the extent and intensity of gender bias in AI-generated text while maintaining the same stereotypical direction as English.

## 5.4 CROSS-LINGUISTIC COMPARISON

In order to find parallels and discrepancies in gender bias tendencies between the two languages, this section compares the English and Italian results in a methodical manner. The comparison looks at stereotype conformity, bias amount.

### 5.4.1 OVERALL BIAS QUANTITY

The amount of substantial gender bias is the most notable distinction between the two languages. Compared to English, Italian exhibits 2.6 times more significant categories (13 out of 20 versus 5 out of 20). The overall bias quantity comparison is summarized in Table 5.11.

Metric	English	Italian
Total Categories Analyzed	20	20
Significant Categories	5 (25.0%)	13 (65.0%)
Male-Biased Categories	1 (5.0%)	7 (35.0%)
Female-Biased Categories	4 (20.0%)	6 (30.0%)
Non-Significant Categories	15 (75.0%)	7 (35.0%)

Table 5.11: Comparison of Bias Quantity Between English and Italian Datasets

There are seven times as many male-biased categories (7 versus 1) and 1.5 times as many female-biased categories (6 versus 4) in the Italian sample. This implies that in the Italian setting, gender-specific language patterns are more noticeable.

### 5.4.2 STEREOTYPE CONFORMITY COMPARISON

English and Italian are the same when it comes to following stereotypes. English follows stereotypes completely one hundred percent of the time. Italian also follows stereotypes completely one hundred percent. You can see how they compare in Table 5.12.

Both languages exhibit perfect stereotype conformity for all trait types. All significant agentic categories are biased toward males as expected by stereotypes, and all significant communal categories are biased toward females as expected. This perfect alignment indicates that when gender bias emerges in

#### 5.4. CROSS-LINGUISTIC COMPARISON

Metric	English	Italian
Conformity Rate	100% (5/5)	100% (13/13)
Conforming Categories	5	13
Counter-Stereotypical Categories	0	0
<b>By Trait Type:</b>		
Agentic Conformity	100% (1/1)	100% (7/7)
Communal Conformity	100% (4/4)	100% (6/6)

Table 5.12: Comparison of Stereotype Conformity Between English and Italian Datasets

either language, it consistently follows traditional gender stereotypes without any counter-stereotypical patterns.

#### 5.4.3 CATEGORY-LEVEL COMPARISON

Patterns of agreement and disagreement between the two languages can be found by comparing them category by category. The full comparison for all 20 categories is shown in Table 5.13.

Category	Type	English	Italian
achievement	Agentic	NS	Male
power	Agentic	NS	Male
anger	Agentic	NS	NS
confusion	Agentic	NS	Male
business	Agentic	NS	NS
work	Agentic	NS	Male
money	Agentic	NS	Male
traveling	Agentic	NS	NS
aggression	Agentic	Male	Male
strength	Agentic	NS	NS
independence	Agentic	NS	NS
pride	Agentic	NS	Male
affection	Communal	Female	Female
party	Communal	Female	Female
friends	Communal	Female	Female
family	Communal	NS	NS
communication	Communal	NS	NS
help	Communal	NS	Female
politeness	Communal	NS	Female
sympathy	Communal	Female	Female

Table 5.13: Category-by-Category Comparison (NS = Not Significant)

This comparison shows a number of trends. First, five categories—aggression

(male-biased), affection (female-biased), party (female-biased), friends (female-biased), and sympathy (female-biased)—are significant in both languages with the same direction. Second, although zero categories are important only in English, six categories—achievement, power, confusion, work, money, pride, assistance, and politeness—are important only in Italian. Third, seven categories—anger, business, travel, strength, independence, family, and communication—are not relevant in either language. Fourth, when significant, all significant categories exhibit the same stereotyped orientation in both languages.

#### 5.4.4 THE AGGRESSION CATEGORY: CONSISTENT PATTERN

In both languages, the category of aggressiveness consistently shows male bias. Aggression is male-biased in both English ( $p = 0.004$ ,  $d = +0.035$ , ratio = 104.3%) and Italian ( $p < 0.001$ ,  $d = +0.066$ , ratio = 131.0%), supporting the stereotype that aggressive language is associated with males.

Since aggression is one of the few agentic categories that shows considerable bias in both languages in the same direction, this consistency is notable. The Italian impact is almost double and as high as the English effect ( $d = +0.066$  versus  $d = +0.035$ ), showing that the male-aggressive stereotype is amplified by grammatical gender marking. Table 5.14 displays specific statistics for this comparison. Table 5.14 provides detailed statistics for this comparison.

Metric	English	Italian
Male Mean Score	0.626	0.074
Female Mean Score	0.600	0.056
Difference	+0.026 (M higher)	+0.018 (M higher)
p-value	0.004	<0.001
Cohen's d	+0.035	+0.066
Bias Ratio	104.3%	131.0%
Direction	Male	Male
Conforms to Stereotype?	Yes	Yes

Table 5.14: Detailed Comparison of Aggression Category Across Languages

This stereotype appears to be stable across linguistic contexts, as evidenced by the consistent male bias in aggression in both languages, with Italian exhibiting a stronger manifestation of the identical pattern seen in English.

## 5.4. CROSS-LINGUISTIC COMPARISON

### 5.4.5 COMMUNAL TRAITS: STRONG CROSS-LINGUISTIC AGREEMENT

Strong cross-linguistic agreement is seen in communal features. In both English and Italian, there are four important communal categories that point in the same direction: attachment, party, friends, and sympathy are all biased toward females. Furthermore, two additional female-biased communal categories—help and politeness—that are not prominent in English are displayed in Italian.

The four overlapping communal categories are detailed in Table 5.15.

Category	English		Italian	
	Cohen's d	p-value	Cohen's d	p-value
affection	-0.045	<0.001	-0.171	<0.001
party	-0.021	0.001	-0.049	0.002
friends	-0.047	<0.001	-0.080	<0.001
sympathy	-0.029	0.010	-0.066	<0.001

Table 5.15: Comparison of Four Overlapping Communal Categories

All four categories exhibit female bias in both languages, with Italian consistently having greater effect sizes. The effect size of affection in Italian ( $d = -0.171$ ) is 3.8 times greater than in English ( $d = -0.045$ ). According to this pattern, communal stereotypes are prevalent in all languages, with Italian exhibiting a more pronounced manifestation of them.

### 5.4.6 EFFECT SIZE COMPARISON

Compared to English, Italian typically exhibits bigger effect sizes. There is a 2.1-fold difference between the mean absolute impact size across significant categories in Italian (0.067) and English (0.032). Compared to the largest effect in English (friends,  $d = -0.047$ ), the largest effect in Italian (affection,  $d = -0.171$ ) is 3.6 times greater.

A visual comparison of impact sizes for all significant categories in both languages is shown in Figure 5.5. Blue bars represent English effect sizes, orange bars represent Italian effect sizes. Positive values indicate male-biased categories, negative values indicate female-biased categories. Italian shows generally larger effect sizes and more significant categories overall.

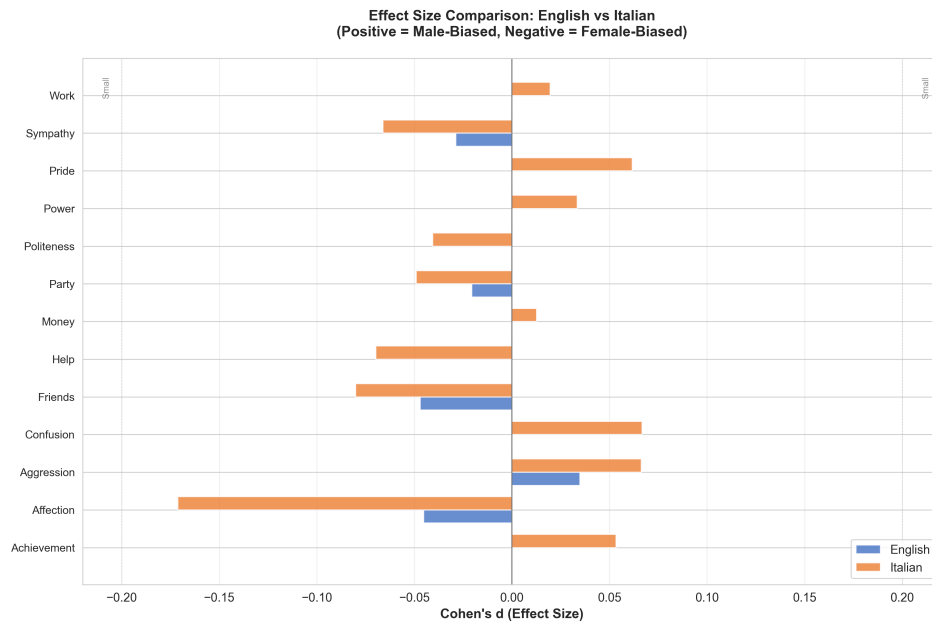


Figure 5.5: Comparison of effect sizes (Cohen's  $d$ ) for significant categories in English and Italian datasets.

### 5.4.7 SUMMARY OF CROSS-LINGUISTIC FINDINGS

The cross-linguistic study indicates significant differences in the number and severity of gender bias between English and Italian, as well as identical patterns in bias direction. Italian has 2.6 times more significant categories (65

The crucial conclusion is that when gender prejudice occurs in either language, it consistently adheres to traditional gender stereotypes: agentic features favor males, while communal traits favor females. Languages differ in the degree and severity of prejudice rather than its direction. Italian displays greater effect sizes and more categories that attain statistical significance.

Communal qualities show excellent cross-linguistic agreement, with four categories (affection, party, friends, and sympathy) exhibiting consistent female bias in both languages. The same pattern is shown in agentic traits: aggression is male-biased in both languages, with Italian exhibiting a greater effect ( $d = +0.066$ ) than English ( $d = +0.035$ ).

These findings indicate that grammatical gender marking in Italian occupational titles increases the prevalence and amount of gender prejudice while retaining the same stereotyped direction as in English. Language-specific factors determine how strongly gender stereotypes appear in AI-generated text, but not their fundamental orientation.

## **5.5** CHAPTER SUMMARY

This chapter presented the empirical data from an experiment examining gender bias in AI-generated interview responses in English and Italian. The main findings are given below.

In the English dataset (29,411 responses), 5 out of 20 categories (25 percent) shows strong gender bias, with perfect stereotype conformance (100 percent). The majority of the prejudice was found in community qualities, which were biased against females in four of the five significant categories. All effect sizes were minor ( $|d| < 0.2$ ), indicating slight differences.

Similar to English, 13 out of 20 categories (65 percent) in the Italian dataset (29,411 responses) exhibited complete stereotype conformance (100 percent). Both agentic and communal features were biased, with six communal categories favoring women and seven agentic categories favoring men. The greatest effect (affection,  $d = -0.171$ ) was much bigger than any English effect, with effect sizes ranging from small to small-medium.

Italian has 2.6 times more gender bias than English in terms of the number of significant categories and the size of effect sizes, according to the cross-linguistic study. Nonetheless, stereotype compliance is 100

These results imply that while gender bias in AI-generated text follows consistent stereotypical patterns across languages, language specific factors, specifically, grammatical gender marking in Italian occupational titles, amplify the degree and magnitude of bias while preserving the same stereotypical direction. Chapter 6 discusses the relevance of these findings for comprehending and reducing gender bias in multilingual AI systems.

# 6

## Discussion

### 6.1 OVERVIEW

The empirical results reported in Chapter 5 are interpreted in this chapter, along with their implications for comprehending gender bias in large language models. Examining gender-differentiated language patterns in DeepSeek-generated interview responses, comparing linguistic patterns in English and Italian, and assessing the practical significance of observed differences are the three research questions that form the basis of the discussion.

The results suggest that both languages have systematic gender-differentiated verbal patterns, with Italian exhibiting significantly greater bias than English. All significant biases conform perfectly to traditional gender stereotypes, with agentic traits consistently associated with male personas and communal traits with female personas. These findings imply that the DeepSeek model reproduces gender preconceptions in professional interview settings after internalizing them from its training data.

### 6.2 INTERPRETATION OF ENGLISH RESULTS (RQ1)

The first research question asked: *How do DeepSeek-generated interview responses exhibit gender bias and align with gender stereotypes?*

According to the examination of the English dataset, DeepSeek produces gender-differentiated language patterns with perfect stereotype compliance in

## 6.2. INTERPRETATION OF ENGLISH RESULTS (RQ1)

25% of stereotype-related criteria (5 out of 20). This result suggests that although the model generates generally fairly balanced replies, systematic biases appear in particular linguistic areas associated with assertiveness and interpersonal connections.

### 6.2.1 PATTERN OF BIAS IN ENGLISH

The English bias pattern is distinguished by an asymmetry in agentic and communal categories. Four communal categories (affection, party, friends, and sympathy) exhibit strong female bias, but just one agentic category (aggression) exhibits significant male prejudice. This 1:4 ratio indicates that the model's gender stereotyping is more noticeable when it comes to interpersonal and emotional language than when it comes to professional competence language.

This tendency is consistent with research in organizational psychology demonstrating that in professional settings, gender stereotypes frequently show up more strongly in assessments of interpersonal skills than technical ability [37, 12]. In most professional domains, the model seems to have learned that male candidates should be portrayed more neutrally, while female candidates should be described using more relationship-oriented and emotionally expressive language.

One notable example of a substantial agentic bias is aggression. Although the model more strongly links aggression to male personas, aggression is not often regarded as a desirable professional trait. This implies that the model has internalized potentially harmful gender connections found in training data in addition to positive professional preconceptions.

### 6.2.2 EFFECT SIZES AND PRACTICAL SIGNIFICANCE

The English dataset's effect sizes are all small ( $|d| < 0.2$ ), suggesting that although changes are statistically significant, they are not dramatic. Friends have the most impact ( $d = -0.047$ ), with affection coming in second ( $d = -0.045$ ).

Compared to what may be anticipated from previous language models, these small impact sizes indicate that DeepSeek's bias in English is comparatively low. However, when responses are produced at scale, minor effects can still have a significant cumulative impact. Even little systematic variations could contribute to the persistence of gender stereotypes in professional communication

if thousands of job applicants employ AI-generated interview responses.

Concerns regarding detection thresholds are also raised by the tiny amount of impacts. More complex analysis methods or lexical frameworks may show additional patterns that the current methodology does not capture.

### 6.2.3 OCCUPATIONAL CONTEXT EFFECTS

According to the job-level analysis, prejudice intensity varies significantly between occupational settings, ranging from 1.7 percent in Education and Training to 11.7 percent in Retail and Service. Gender bias in AI-generated text varies systematically by professional domain, as evidenced by this 6.9-fold difference.

Service-oriented occupations (Retail and Service, Food and Hospitality) have the greatest prejudice, whereas professional positions requiring specific credentials (Education and Training, Healthcare and Medical) have the least bias. This pattern implies that descriptions of customer-facing and service roles in the model's training data may have more gender stereotypes than descriptions of professional roles with credentials.

It's interesting to note that bias severity and occupational gender dominance do not systematically correlate. Male-dominated and female-dominated categories appear throughout the distribution, showing that bias arises from culturally established stereotyped associations rather than basic demographic representation in training data.

## 6.3 INTERPRETATION OF ITALIAN RESULTS (RQ2)

The second research question asked: *Are there differences in the frequency and directional distribution of gender-differentiated language between English and Italian?*

The Italian dataset has a far higher gender bias than the English dataset, with 65% of categories exhibiting significant differences compared to 25% in English. One of the study's most startling results is this 2.6-fold increase, which implies that language-specific factors significantly boost gender bias in AI-generated writing.

### 6.3. INTERPRETATION OF ITALIAN RESULTS (RQ2)

#### 6.3.1 AMPLIFICATION OF BIAS IN ITALIAN

The Italian results reveal considerable bias in 13 of 20 areas, including 7 agentic categories (achievement, power, confusion, labor, money, aggression, pride) and 6 communal categories (affection, party, friends, aid, politeness, sympathy). This reflects a far more extensive pattern of gender difference than that found in English.

Several variables could lead to this amplification. First, Italian has grammatical gender specification on occupational titles, making gender more linguistically prominent throughout the prompt and response. When the model generates the prompts with gendered job titles, such as *ingegnere* (male engineer) vs *ingegnera* (female engineer), the grammatical gender may prime gender-associated linguistic patterns more strongly than in English.

Second, Italian training data may contain more gender stereotypes than English data. If there is a greater gender differentiation in how male and female professionals are described in Italian-language professional discourse, the model will learn and imitate these patterns.

Third, the translation process from Italian to English for lexical analysis may have contributed some artifacts, but the consistent direction of all biases implies that this is not the key reason.

#### 6.3.2 LARGER EFFECT SIZES IN ITALIAN

Italian has large effect sizes than English, with an average absolute effect size of 0.067 against 0.032 in English. The highest effect in Italian (affection,  $d = -0.171$ ) is 3.8 times that in English (friends,  $d = -0.047$ ).

These greater effect sizes show that gender distinction in Italian is not only more common, but also more striking when it does occur. The affection category has a very significant female bias in Italian, implying that the model correlates female personas with emotionally expressive language far more strongly in Italian than in English.

The bigger effect sizes in Italian shift various biases from the "small" to "small-medium" range, implying that the practical importance of bias may be stronger in Italian-language implementations of the model.

### 6.3.3 CROSS-LINGUISTIC CONSISTENCY

Despite the size changes, the direction of bias remains relatively similar across languages. All five categories with considerable prejudice in English have significant bias in the same direction in Italian. This perfect directional agreement implies that the underlying gender preconceptions are culturally universal rather than language particular.

The consistency is particularly noticeable in the hostility category, which exhibits masculine bias in both languages with identical patterns. This shows that the link between masculine personalities and aggressive language persists across linguistic situations.

The four communal categories that bridge languages (affection, party, friends, and sympathy) all have a female bias in both English and Italian, demonstrating that the link between female personalities and interpersonal/emotional language is consistent across the two languages analyzed.

### 6.3.4 ITALIAN-SPECIFIC BIASES

Italian has a considerable bias in eight categories that are not prominent in English: achievement, authority, confusion, effort, money, pride (agentic), and help, politeness (community). These Italian-specific biases indicate that grammatical gender marking may enhance stereotypes in a broader variety of professional and interpersonal areas.

The appearance of major bias in performance, work, and power in Italian is especially notable, given that these are essential professional competency areas. This implies that the Italian model distinguishes between male and female personalities not only in interpersonal areas, but also in descriptions of professional performance and authority.

## 6.4 STEREOTYPE CONFORMITY ANALYSIS

Both languages exhibit perfect stereotype conformity (100 percent), which means that all significant biases correspond to standard gender preconceptions. Agentic categories are constantly skewed toward men, while communal categories are continuously biased toward women, with no counter-stereotypical patterns found.

## 6.5. COMPARISON WITH PREVIOUS RESEARCH

The flawless conformity is both stunning and frightening. It shows that the DeepSeek model generates gender-differentiated language in ways that reinforce rather than question existing gender roles. From its training data, the model learnt statistical connections between gender and professional/interpersonal qualities, which it systematically reproduces.

The lack of any counter-stereotyped biases shows that the model's training data is insufficiently rich in examples of gender-atypical professional descriptions to overcome the dominating stereotypical patterns. This finding highlights the challenge of achieving fairness in language models trained on real-world text that reflects existing social biases.

## **6.5** COMPARISON WITH PREVIOUS RESEARCH

The results of this research may be compared to a number of related studies that analyze the problem of gender bias in large language models. The closest point of comparison is the original study by Kong et al. (2024), to which the current thesis is the replica and extension of its experimental framework. The two other recent studies, namely Kalhor and Bahrak (2025) and Du et al. (2025) provide more context by analyzing gender bias in the DeepSeek in various linguistic and professional contexts.

Kong et al. (2024) discovered that GPT-3.5, GPT-4, and Claude show the same gender-differentiated patterns in interview responses, where male personalities are related to agentic language and female ones to communal language. The current paper demonstrates that the same trends can be observed in DeepSeek, indicating that gender biasness in generating response to interview questions is not a specific property of particular architecture or training methods but a widespread feature of large language models. Although direct numerical comparison is constrained by the application of various lexical analysis tools (LIWC versus Empath), the directional similarity in the results of models is significant.

Kalhor and Bahrak (2025) [51] investigated gender stereotypes in four multilingual LLMs with DeepSeek R1 in both Persian and English and discovered that all models had gender stereotypes but with more differences in Persian than in English. Persian, as Italian, is a grammatically gendered language, and their observation that grammatical gender enhances bias is also in line with the

cross-linguistic trend found in the current research.

Du et al. (2025) [52] benchmarked gender bias between six LLMs, such as DeepSeek-R1 in an educational feedback setting. They observed gender differences in their semantic responses to gender cues, that is, more autonomy-supportive language in male cues and more controlling language in female cues. This agentic-male and communal-female pattern is in line with the stereotypical patterns observed in the current research which indicates that the gender-differentiated language generation of DeepSeek is not limited to the sphere of job interviews but also spread to other areas of professional activities.

Collectively, these comparisons indicate that gender bias in DeepSeek is an organized and strong effect that is present in various languages, fields, and methodological systems.

## **6.6** THEORETICAL IMPLICATIONS

The findings of this study have several important theoretical implications for understanding bias in large language models.

### **6.6.1** SOCIAL ROLE THEORY AND LANGUAGE MODELS

The discovered biases perfectly agree with predictions from social role theory [24], implying that language models acquire and reproduce culturally entrenched gender preconceptions. The persistent connection of agentic features with male personas and communal traits with female personas corresponds to trends found in decades of social psychology research on gender stereotypes.

This finding lends support to the notion that language models act as "statistical mirrors" of their training data, mirroring the social patterns and biases found in human-generated text. The models do not generate new stereotypes; instead they learn and reinforce current ones.

### **6.6.2** GRAMMATICAL GENDER AND BIAS AMPLIFICATION

The significantly higher bias in Italian compared to English lends empirical support to the theory that grammatical gender marking increases gender stereotyping in language model outputs. When gender is marked grammatically

## 6.7. PRACTICAL IMPLICATIONS

throughout prompts and responses, it appears to strengthen gender-associated language patterns.

This finding has important implications for understanding how linguistic structure interacts with social cognition in neural language models. It implies that models trained on grammatically gendered languages may require bias mitigation measures other than those used for English.

### **6.6.3** OCCUPATIONAL STEREOTYPES IN AI SYSTEMS

The variance in bias intensity across occupational categories shows that language models have learned both generic gender stereotypes and occupation-specific stereotyped connections. The model's training data exhibits distinct patterns of gender-differentiated language across professional domains, as evidenced by the larger bias in service-oriented jobs relative to credentialed professional positions.

This result implies that rather than being implemented consistently across all professional domains, bias mitigation techniques may need to be customized for certain occupational situations.

## **6.7** PRACTICAL IMPLICATIONS

The findings of this study have important practical implications for the deployment and use of large language models in professional contexts.

### **6.7.1** IMPLICATIONS FOR AI-ASSISTED JOB APPLICATIONS

The use of AI writing tools in job application situations will be immediately impacted by the discovery that DeepSeek generates gender-specific interview responses. The language used in cover letter and interview responses are going to appear in the language of job seekers that activate gender stereotypes for employers.

Female applicants who seek AI help might be answered more affectively, depending on the trends that were observed in this study.

Even if individual application receivers are oblivious to the biases within recommendations provided by AI, these minor variations might accumulate

over multiple job applications and therefore contribute to the lasting of gendered stereotypes in written workplace interactions.

### **6.7.2** IMPLICATIONS FOR MULTILINGUAL AI SYSTEMS

The finding that Italian has dramatically higher bias than English has several implications for using language models in multilingual environments. Any organization utilizing NLP writing tools in various languages should be aware that there exists a wide range of bias levels, and statistically gendered languages may have higher levels of stereotypes.

Indicates that the bias auditing and mitigation methods found for one language may not quite adapt well to other languages. This raises the need for language dependent bias auditing and mitigation approaches for multilingual systems.

### **6.7.3** IMPLICATIONS FOR MODEL DEVELOPMENT

The fact that pronouns and occupations were stereotyped equally well in both languages is a further indicator that current training and receptor alignment techniques are inadequate for eliminating the appearance of gender stereotypes in generated text. Model developers might have to take more specific measures to avoid stereotypical gender-profession associations.

Possible approaches could include:

- Enlarging the training data with counterstereotypical examples
- Applying the bias-aware fine-tuning techniques
- Devising prompts that are gender insensitive
- Explicit fairness constraints (e.g., we require the model to all using- or to avoid all in this sentence) during generation

Given the greater bias in Italian, models trained on grammatically gendered languages could need more bias mitigation techniques than those created for English.

## **6.8** LIMITATIONS

Several limitations of this study should be acknowledged when interpreting the findings.

## 6.8. LIMITATIONS

### 6.8.1 METHODOLOGICAL LIMITATIONS

Firstly, only one language model has been used (DeepSeek), so conclusions cannot be generalized to other language models. There were similar trends to previous work with other models, but would need to run a similar experiment on a different language model to compare.

Second, the lexical analysis is based on the Empath framework, which may capture distinct features of language than techniques such as LIWC. While Empath provides a transparent and reproducible approach, the precise categories and criteria used may determine which biases are found.

Third, the Italian responses were translated into English for lexical analysis, which may have introduced artifacts or obscured some language-specific nuances. While the uniform direction of all biases indicates that translation is not the key driver of results,

### 6.8.2 SCOPE LIMITATIONS

The evidence and analysis in this thesis solely pertains to gender bias and does not take account of other types of bias. Race, ethnicity, age and disability may come into play in interactions with gender bias.

This research is of two languages (English and Italian). The distribution of the phenomenon may vary in non-European languages and languages with more or less grammatical genders.

There is also the issue of artificiality of the environment. The fake interview may not be how the real user would go about using the AI in the real world, and the user may have changed the model suggestions.

### 6.8.3 INTERPRETATION LIMITATIONS

This study highlights statistical regularities in model outputs, but doesn't provide any way of inferring how these regularities emerge. The tendency for training data bias stands out as the most plausible explanation, but the way in which models' architecture or fine-tuning method operate could also be factors.

The practical significance of small effect sizes cannot be evaluated. Again there are statistically significant differences, but how this translates into actual hiring decisions or how it affects a person's view of a job?

This study does not investigate whether the biases found have any relevance to actual human perception. To do this, an experiment would have to be carried out on human raters to see if the variations in language identified by computational analysis are detectable and important in real professional situations.

## 6.9 SUMMARY

This discussion has interpreted the empirical data in light of the three study topics, explored theoretical and practical consequences, and acknowledged significant limits. The major insights include:

- DeepSeek develops gender-differentiated verbal patterns that perfectly match traditional gender preconceptions.
- Italian has 2.6 times more bias than English, indicating that grammatical gender increases stereotyping.
- Bias is systematic across occupational contexts, with service roles displaying the biggest bias.
- Effect sizes are normally small, but they can have a cumulative impact at scale.
- Existing training methods appear insufficient to minimize gender bias in professional language creation.

These findings help us understand how large language models repeat societal biases and highlight the need for more study into bias identification and reduction, particularly in multilingual environments.





# Conclusions

## 7.1 OVERVIEW

This thesis explored gender bias in interview responses generated by the DeepSeek large language model, focusing on linguistic patterns in English and Italian. The study used a controlled experimental design, with 29,411 interview responses per language across 60 occupations and 5 standardized interview questions.

The important discovery is that DeepSeek has systematic gender-differentiated patterns that perfectly match traditional gender preconceptions, with Italian displaying 2.6 times greater bias than English.

## 7.2 SUMMARY OF KEY FINDINGS

The English dataset revealed gender-differentiated patterns in 25% of categories (5 of 20), with 100% stereotype conformance. The Italian dataset had 65% significant categories (13 out of 20), along with full stereotype conformance. All major biases are consistent with traditional gender expectations: agentic features for men, communal traits for women.

## 7.3 RESEARCH QUESTIONS ANSWERED

**RQ1:** DeepSeek generates systematic gender-differentiated linguistic patterns that align perfectly with traditional gender stereotypes, achieving 100% stereotype conformance in both English and Italian. The agentic and communal categories are always biased on male and female personas respectively. At the job level, bias intensity varies substantially across occupational contexts, with service-oriented roles showing the highest bias and credentialed professional roles showing the lowest.

**RQ2:** Italian has 2.6 times more bias frequency than English, yet the directional constancy is perfect across languages.

**RQ3:** Effect sizes are small but systematic, with practical implications varying according to deployment scale and language.

## 7.4 CONTRIBUTIONS

This thesis makes the following contributions: (1) one of the first systematic evaluations of DeepSeek for gender bias in the context of job interview responses and work selection, to the best of the author’s knowledge, (2) new evidence that grammatical gender amplifies bias, (3) the demonstration of Empath as an open-source alternative for bias evaluation, and (4) a thorough cross-linguistic methodology comparing English and Italian.

## 7.5 LIMITATIONS

Limitations include the use of a single model, just two languages, Empath framework specialization, translation effects, a gender-exclusive focus, lexical analysis, and a simulated setting. Furthermore, the use of Empath rather than LIWC (used in the original study by Kong et al., 2024) means that direct numerical comparisons with the original findings are not possible, as the two frameworks use different category definitions and lexical coverage.

The study is limited to the DeepSeek model only, and the findings may not generalize to other large language models such as GPT-4 or Claude. Only two languages were examined; the patterns observed may differ in other grammatically gendered languages such as Spanish, French, or German. The lexical

analysis relies entirely on the Empath framework, which may not capture all relevant linguistic dimensions of gender bias. The Italian responses were translated to English before analysis, which may have introduced translation artifacts. The study focuses exclusively on binary gender and does not address intersectional biases involving race, age, or disability. Finally, the simulated interview setting may not fully reflect real-world AI-assisted job application scenarios.

## 7.6 FUTURE WORK

Future study should look into other models, more language pairs, intersectional bias, temporal dynamics, mitigation measures, human perception studies, and different analysis methodologies.

## 7.7 CLOSING REMARKS

This work demonstrates that current LLMs perpetuate gender stereotypes even in professional contexts, and that marking grammatical gender intensifies bias. Because the perfect stereotype conformation for both languages, it implies that contemporary training strategies do not address ingrained gender bias. To produce fair multilingual AI systems, language-specific bias audit and reduction strategy are required.

It is imperative that we continue to research in the area of AI so that tools don't get in the way of gender equality in the workplace as they are more and more incorporated into working life.

## 7.8 USE OF GENERATIVE AI TOOLS

This thesis involved the use of AI tools at different stages of the work. The focus of this research was DeepSeek (deepseek-chat model) — all of the simulated interview responses used in this study were produced using the DeepSeek API. Translating Italian interview responses to English was also done using the same model during the data preprocessing stage.

The writing process involved the use of ChatGPT to facilitate language editing and  $\text{\LaTeX}$  formatting. It was not used to produce research findings, conduct

## 7.9. DATA AND CODE AVAILABILITY

analysis, or draw conclusions. All experimental design, data collection, statistical testing, and interpretation of results were carried out entirely by the author.

The author reviewed all AI-assisted content and takes full responsibility for the accuracy and integrity of this thesis.

### **7.9** DATA AND CODE AVAILABILITY

To promote transparency and reproducibility, all materials related to this thesis are publicly accessible. The experimental data, including the English and Italian interview responses generated by the DeepSeek model, and the full code of data preprocessing, empirical lexical analysis with the Empath framework, statistical testing, and creation of figures are stored on GitHub at:

[https://github.com/memonhufsa205/Thesis\\_Dataset](https://github.com/memonhufsa205/Thesis_Dataset)

## References

- [1] Haein Kong et al. "Gender Bias in LLM-generated Interview Responses". In: *arXiv preprint arXiv:2410.20739* (2024). URL: <https://arxiv.org/abs/2410.20739>.
- [2] Wayne Xin Zhao et al. "A survey of large language models". In: *arXiv preprint arXiv:2303.18223* (2023). URL: <https://arxiv.org/abs/2303.18223>.
- [3] Tom B. Brown et al. "Language Models are Few-Shot Learners". In: *Advances in Neural Information Processing Systems*. 2020.
- [4] Saurabh Bhausahab Zinjad et al. "ResumeFlow: An LLM-facilitated pipeline for personalized resume generation and refinement". In: *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2024, pp. 2781–2785.
- [5] Emily M Bender et al. "On the dangers of stochastic parrots: Can language models be too big?" In: *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 2021, pp. 610–623.
- [6] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. "Semantics Derived Automatically from Language Corpora Contain Human-like Biases". In: *Science* 356.6334 (2017).
- [7] Yixin Wan et al. "'Kelly is a warm person, Joseph is a role model': Gender biases in LLM-generated reference letters". In: *arXiv preprint arXiv:2310.09219* (2023). URL: <https://arxiv.org/abs/2310.09219>.
- [8] Juan M Madera, Michelle R Hebl, and Randi C Martin. "Gender and letters of recommendation for academia: agentic and communal differences". In: *Journal of Applied Psychology* 94.6 (2009), p. 1591.

## REFERENCES

- [9] Haozhe An et al. “Do large language models discriminate in hiring decisions on the basis of race, ethnicity, and gender?” In: *Findings of the Association for Computational Linguistics: ACL 2024*. 2024. URL: <https://arxiv.org/abs/2406.10486>.
- [10] Abel Salinas et al. “The unequal opportunities of large language models: Examining demographic biases in job recommendations by chatgpt and llama”. In: *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*. 2023, pp. 1–15.
- [11] Sabine Sczesny, Christa Nater, and Alice H Eagly. “Agency and communion: Their implications for gender stereotypes and gender identities”. In: *Agency and communion in social psychology*. Routledge, 2018, pp. 103–116.
- [12] Laurie A Rudman and Peter Glick. “Prescriptive gender stereotypes and backlash toward agentic women”. In: *Journal of social issues* 57.4 (2001), pp. 743–762.
- [13] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. MIT Press, 2019.
- [14] N. Mehrabi et al. “A Survey on Bias and Fairness in Machine Learning”. In: *ACM Computing Surveys* 54.6 (2021), pp. 1–35. DOI: 10.1145/3457607.
- [15] Ioana M Latu, Marianne Schmid Mast, and Tracie L Stewart. “Gender biases in (inter) action: The role of interviewers’ and applicants’ implicit and explicit stereotypes in predicting women’s job interview outcomes”. In: *Psychology of Women Quarterly* 39.4 (2015), pp. 539–552.
- [16] Sharon L Segrest Purkiss et al. “Implicit sources of bias in employment interview judgments and decisions”. In: *Organizational Behavior and Human Decision Processes* 101.2 (2006), pp. 152–167.
- [17] Matthew Hutson. “Artificial Intelligence Faces Reproducibility Crisis”. In: *Science* 359.6377 (2018), p. 725. DOI: 10.1126/science.359.6377.725.
- [18] Ryan L Boyd et al. *The development and psychometric properties of LIWC-22*. Tech. rep. University of Texas at Austin, 2022.
- [19] Erik Fast, Binbin Chen, and Michael S. Bernstein. “Empath: Understanding Topic Signals in Large-Scale Text”. In: *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. San Jose, California, USA: ACM, 2016, pp. 4647–4657.

- [20] Ashish Vaswani et al. “Attention is All You Need”. In: *Advances in Neural Information Processing Systems*. 2017.
- [21] Alec Radford et al. *Improving Language Understanding by Generative Pre-Training*. Tech. rep. OpenAI, 2018.
- [22] Alec Radford et al. *Language Models are Unsupervised Multitask Learners*. Tech. rep. OpenAI, 2019.
- [23] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *NAACL-HLT*. 2019.
- [24] Alice H. Eagly and Valerie J. Steffen. *Sex Differences in Social Behavior: A Social-Role Interpretation*. Lawrence Erlbaum Associates, 1987.
- [25] Henry B Mann and Donald R Whitney. “On a test of whether one of two random variables is stochastically larger than the other”. In: *The Annals of Mathematical Statistics* 18.1 (1947), pp. 50–60.
- [26] Jacob Cohen. *Statistical Power Analysis for the Behavioral Sciences*. 2nd. Lawrence Erlbaum Associates, 1988.
- [27] Tolga Bolukbasi et al. “Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings”. In: *Advances in Neural Information Processing Systems*. 2016.
- [28] Tomas Mikolov et al. “Efficient Estimation of Word Representations in Vector Space”. In: *ICLR Workshop*. 2013.
- [29] Jeffrey Pennington, Richard Socher, and Christopher Manning. “GloVe: Global Vectors for Word Representation”. In: *EMNLP*. 2014.
- [30] Rachel Rudinger et al. “Gender bias in coreference resolution”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. 2018, pp. 8–14.
- [31] Jieyu Zhao et al. “Gender bias in coreference resolution: Evaluation and debiasing methods”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. 2018, pp. 15–20.
- [32] Keita Kurita et al. “Measuring Bias in Contextualized Word Representations”. In: *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. 2019.

## REFERENCES

- [33] Emily Sheng et al. “The Woman Worked as a Babysitter: On Biases in Language Generation”. In: *EMNLP*. 2019.
- [34] Moin Nadeem, Anna Bethke, and Siva Reddy. “StereoSet: Measuring stereotypical bias in pretrained language models”. In: *ACL-IJCNLP*. 2021.
- [35] Li Lucy and David Bamman. “Gender and Representation Bias in GPT-3 Generated Stories”. In: *Proceedings of the Third Workshop on Narrative Understanding*. 2021.
- [36] Hadas Kotek, Rikker Dockum, and David Sun. “Gender bias and stereotypes in Large Language Models”. In: *Proceedings of The ACM Collective Intelligence Conference*. 2023.
- [37] Madeline E. Heilman. “Gender stereotypes and workplace bias”. In: *Research in Organizational Behavior* 32.1 (2012), pp. 113–135.
- [38] Corinne A. Moss-Racusin et al. “Science faculty’s subtle gender biases favor male students”. In: *Proceedings of the National Academy of Sciences* 109.41 (2012), pp. 16474–16479.
- [39] Jessica C Babal et al. “Linguistic analysis of pediatric residency personal statements: gender differences”. In: *BMC medical education* 19.1 (2019), pp. 1–9.
- [40] Alysen Demzik et al. “Gender-based differences in urology residency applicant personal statements”. In: *Urology* 150.1 (2021), pp. 2–8.
- [41] Anna Sahlstrom Drury et al. “Gender bias in postgraduate year one pharmacy letters of recommendation”. In: *Journal of the American College of Clinical Pharmacy* 6.5 (2023), pp. 446–457.
- [42] Erryn Tappy et al. “Linguistic differences by gender in letters of recommendation for minimally invasive gynecologic surgery fellowship applicants”. In: *Journal of Surgical Education* 79.4 (2022), pp. 928–934.
- [43] Deanna M Kaplan et al. “What’s in a name? experimental evidence of gender bias in recommendation letters generated by chatgpt”. In: *Journal of Medical Internet Research* 26.1 (2024), e51837.
- [44] Khanh Cao Nguyen et al. “Rethinking skill extraction in the job market domain using large language models”. In: *arXiv preprint arXiv:2402.03832* (2024). URL: <https://arxiv.org/abs/2402.03832>.
- [45] OpenAI. *GPT-4 Technical Report*. Tech. rep. OpenAI, 2023.

- [46] Su Lin Blodgett et al. "Language (Technology) is Power: A Critical Survey of "Bias" in NLP". In: *ACL*. 2020.
- [47] Gavin Abercrombie and Dirk Hovy. "What can we do about it? Mitigating gender bias in NLP". In: *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. 2019.
- [48] Odd Erik Gundersen and Sigbjørn Kjensmo. "State of the Art: Reproducibility in Artificial Intelligence". In: *AAAI*. 2018.
- [49] Joelle Pineau et al. "Improving Reproducibility in Machine Learning Research". In: *Journal of Machine Learning Research*. 2021.
- [50] Social Security Administration. *Top names over the last 100 years*. <https://www.ssa.gov/oact/babynames/decades/century.html>. Accessed: 2024-07-15. 2024.
- [51] Ghazal Kalhor and Behnam Bahrak. "Probing Gender Bias in Multilingual LLMs: A Case Study of Stereotypes in Persian". In: *Proceedings of the 9th Widening NLP Workshop*. Suzhou, China: Association for Computational Linguistics, 2025, pp. 19–27. DOI: 10.18653/v1/2025.winlp-main.3. URL: <https://aclanthology.org/2025.winlp-main.3>.
- [52] Yishan Du et al. "Benchmarking Educational LLMs with Analytics: A Case Study on Gender Bias in Feedback". In: *arXiv preprint arXiv:2511.08225* (2025). URL: <https://arxiv.org/abs/2511.08225>.



# Acknowledgments

Above all, I am deeply grateful to Allah (SWT) for His endless mercy, guidance, and blessings, which have granted me the strength and determination to accomplish this work. Every step of this journey has been possible through His will.

I would like to express my sincere appreciation to my supervisor, Prof. Antonio Roda, who was always supportive, patient, and insightful in guiding me during the entire course of this research. His experience and insightful comments have been instrumental in polishing my ideas and completing this thesis.

My heartfelt thanks go to the University of Padova and the Department of Information Engineering, which offered me a dynamic learning environment, resources, and facilities required to complete this study.

My family deserves a special mention; I will always be thankful to the unconditional love and support of my family. I would like to acknowledge them because they believed in me and supported me all my life.

I am also thankful to my friends, who have always been by my side in any circumstances, and have always encouraged, understood, and loved me.

Lastly, I extend my gratitude to all individuals who have contributed, directly or indirectly, to the completion of this thesis. Your support and kindness are sincerely appreciated.

Thank you all.

*Hufsa Memon*  
Padova, 14 April 2026