



**POLITECNICO  
DI TORINO**

**POLITECNICO DI TORINO**

Master Degree course in Engineering and Management

Master Degree Thesis

# **Gender Bias in Generative AI: An Analysis of Recruitment Processes**

## **Supervisors**

Prof. Riccardo COPPOLA

Prof. Marco RONDINA

## **Candidate**

Martina ULLASCI

ACADEMIC YEAR 2024-2025

*In un mondo che ci vuole piccole,  
occupare spazio è un atto rivoluzionario*

– Rebecca Momoli

*A mia mamma*

# Acknowledgements

My first and deepest thanks go to Riccardo Coppola, my supervisor, and Marco Rondina, my co-supervisor, for giving me the opportunity to conduct this research on gender disparity, a topic that is deeply close to my heart and one I care and fight for every day. Without them this work would not have been possible. Thank you for your trust, your availability and for guiding me through a research journey that has given me so much, both professionally and personally.

My heartfelt thanks also go to my mum, my dad and Cekki, my family - biological and chosen - who have been by my side throughout these long years at the Politecnico, through both the good and the difficult times. You are and always will be my safe harbor.

Thanks to all the people with whom I shared one of the most meaningful and beautiful experiences of my life, the occupation of Politecnico, inspired by the university movements for a **Free Palestine**. We spent wonderful months together, showing the world that an alternative to the capitalist and militarized model is possible, together.

Thank you to my long-time friends and to the new ones, in particular to Paciottoli. Each of you has been there in your own way during these years and I love you all.

Finally, thank you to Cavourlandia, my family in Turin. You taught me that it's possible to feel at home, even when you're far from home.

## Abstract

In recent years, generative artificial intelligence (Gen AI) systems have assumed increasingly crucial roles in selection processes, personnel recruitment and analysis of candidates' profiles. However, the employment of large language models risks reproducing, and in some cases amplifying, gender stereotypes and bias already present in labour market. This research aims to evaluate and measure this phenomenon, analysing how a state-of-the-art generative model (GPT-5) suggests occupations and represents the *ideal candidate* based on gender and work experience background, focusing on under 35 years old Italian graduates. The study is organised into two complementary phases. In Phase I, the model has been trained to provide job suggestions to 24 simulated candidates profiles of female and male genders, balanced in age, experience and professional field. The output variables - as job title, industry and descriptive adjectives - were coded using open coding and tested statistically with  $\chi^2$  test. Results show that, although no significant differences emerged in job titles and industry, gendered linguistic patterns exist in the adjectives attributed to female and male candidates, indicating a tendency of the model in associating women to emotional and empathetic traits, while men to strategic and analytic ones. Phase II employed 114 LinkedIn job advertisements, used as prompts for generating textual and visual representations of *ideal candidates*. The analysis of the outputs highlighted a clear gender polarisation: the model assigned 71% of profiles to male and 29% to female gender. The strongest association emerge in technology and engineering sectors, where male candidates prevail and in HR and commercial functions, where female representation prevails. The visual analysis confirms the perpetuation of gender stereotypes: female profiles are more frequently depicted smiling, in approachable postures and dressed elegantly, while men portraits result more focused and assertive with formal clothing style. These results prove that Gen AI models do not simply reflect the gender biases of the training data, but they also amplify them, as the experiment in the labour market context clearly shows. The research raises an ethical question regarding the use of these models, highlighting the need for transparency and bias mitigation strategies to ensure fairness and inclusive representation.



# Contents

<b>List of Figures</b>	4
<b>List of Tables</b>	5
<b>1 Introduction</b>	7
1.1 Context and relevance of generative artificial intelligence . . . . .	7
1.2 From technological design to social inequality . . . . .	8
1.3 The importance of generative AI and work platforms . . . . .	9
1.4 Contribution and purpose of the thesis . . . . .	10
<b>2 Background</b>	13
2.1 Gender bias and social roots of patriarchy . . . . .	13
2.2 Generative AI and gender bias . . . . .	14
2.2.1 Definition of transmission and mitigation of bias . . . . .	15
2.2.2 Definition of types of AI technology . . . . .	16
2.2.3 Transmission and mitigation by technology . . . . .	19
2.2.4 Case study: representation of O*NET occupations . . . . .	22
2.3 Gender bias in recruiting platforms . . . . .	24
<b>3 Methodology</b>	27
3.1 Objectives and research questions . . . . .	27
3.2 Tools and model used . . . . .	27
3.2.1 Text generation model . . . . .	28
3.2.2 Visual generation model . . . . .	28
3.2.3 Primary material sources . . . . .	28
3.2.4 Data collection software . . . . .	28
3.3 Phase I: job suggestions based on simulated candidate profiles . . . . .	28
3.3.1 Construction of the candidate population . . . . .	28
3.3.2 Prompting protocol and output generation . . . . .	32
3.3.3 Data processing and analytical strategy . . . . .	32
3.3.4 Analytical techniques . . . . .	35
3.3.5 Visual summary of Phase I methodology . . . . .	36
3.4 Phase II: ideal candidate generation from real job advertisements . . . . .	38
3.4.1 Construction of the Job advertisements population . . . . .	38
3.4.2 Prompting protocol and output generation . . . . .	40

3.4.3	Data processing and analytical strategy . . . . .	43
3.4.4	Analytical techniques . . . . .	46
3.4.5	Visual summary of Phase II methodology . . . . .	48
<b>4</b>	<b>Results</b>	<b>51</b>
4.1	Results for Phase I: Candidate-driven experiment . . . . .	51
4.1.1	Textual analysis: Job title suggestions . . . . .	51
4.1.2	Textual analysis: Industry suggestions . . . . .	54
4.1.3	Textual analysis: Adjectives suggestions . . . . .	55
4.2	Results for Phase II: Employer-driven experiment . . . . .	58
4.2.1	Overall Gender assignment analysis . . . . .	58
4.2.2	Textual analysis: Job title influence on Gender assignment . . . . .	60
4.2.3	Textual analysis: Industry influence on Gender assignment . . . . .	62
4.2.4	Textual analysis: Adjectives assigned to candidates . . . . .	63
4.2.5	Visual analysis: Posture . . . . .	66
4.2.6	Visual analysis: Facial expression . . . . .	69
4.2.7	Visual analysis: Smiling presence . . . . .	71
4.2.8	Visual analysis: Clothing style . . . . .	73
<b>5</b>	<b>Conclusion</b>	<b>79</b>
5.1	Summary of key findings . . . . .	79
5.2	Critical interpretation . . . . .	80
5.3	Limits of the study and future perspectives . . . . .	80
5.4	Final conclusion . . . . .	80
<b>A</b>	<b>Input materials used for prompt construction</b>	<b>83</b>
A.1	Profiles used as input in Phase I . . . . .	83
A.2	Job advertisements used as input in Phase II . . . . .	83
<b>B</b>	<b>Open coding Tables</b>	<b>87</b>
B.1	Suggested Job title (Phase I) . . . . .	88
B.2	Suggested Industry (Phase I) . . . . .	89
B.3	Suggested Adjectives (Phase I) . . . . .	90
B.4	Suggested Adjectives (Phase II) . . . . .	91
B.5	Suggested Posture adjectives (Phase II) . . . . .	92
B.6	Suggested Facial expression adjectives (Phase II) . . . . .	93
B.7	Suggested Clothing style adjectives (Phase II) . . . . .	94
	<b>Bibliography</b>	<b>95</b>

# List of Figures

3.1	Research methodology flow chart (Phase I). . . . .	37
3.2	Research methodology flow chart (Phase II). . . . .	49
4.1	Distribution of suggested Job title classes by Gender (Phase I). . . . .	52
4.2	Distribution of suggested Industry classes by Gender (Phase I). . . . .	55
4.3	Distribution of suggested Adjective classes by Gender (Phase I). . . . .	56
4.4	Overall distribution of AI-assigned Gender for ideal candidates (Phase II)	59
4.5	Normalized distribution of Gender by Job title classes (Phase II). . . . .	61
4.6	Normalized distribution of Gender by Industry classes (Phase II). . . . .	63
4.7	Normalized distribution of suggested Adjective classes by Gender (Phase II). . . . .	65
4.8	Normalized distribution of suggested Posture classes (Phase II). . . . .	67
4.9	Normalized distribution of suggested Facial expression classes (Phase II).	69
4.10	Normalized distribution of Smiling presence - Gender (Phase II). . . . .	72
4.11	AI-generated portraits of the three ideal candidates for Talent Acquisition Specialist Senior position. . . . .	73
4.12	AI-generated portraits of the three ideal candidates for Product Analyst Senior position. . . . .	73
4.13	Normalized distribution of suggested Facial expression classes (Phase II).	74
4.14	Comparison between female and male candidates portraits . . . . .	76

# List of Tables

2.1	Percentages of representations of women and men according to the model.	23
3.1	Classification of ISCO-08 occupational groups into macro-areas used for simulation. . . . .	30
3.2	Combinations of independent variables within the population. . . . .	31
3.3	Example of prompt input and model output Phase I. . . . .	33
3.4	The combinations of independent variables within the population. . . . .	40
3.5	Example of prompt input and model output Phase II. . . . .	42
4.1	Observed frequencies of job titles for each gender (Phase I). . . . .	53
4.2	$\chi^2$ independence test Job title - Gender (Phase I). . . . .	53
4.3	Final inferential results for Job title (Phase I). . . . .	54
4.4	Observed frequencies Industry - Gender (Phase I). . . . .	54
4.5	$\chi^2$ independence test Industry - (Phase I). . . . .	55
4.6	Final inferential results for Industry suggestions (Phase I). . . . .	56
4.7	Observed frequencies Adjectives - Gender (Phase I). . . . .	56
4.8	$\chi^2$ independence test Adjectives - Gender (Phase I). . . . .	57
4.9	Final inferential results for Adjectives suggestions (Phase I). . . . .	57
4.10	$\chi^2$ goodness-of-fit test Gender distribution (Phase II). . . . .	59
4.11	Final inferential results for Gender distribution (Phase II). . . . .	60
4.12	Observed frequencies Job title - Gender (Phase II). . . . .	60
4.13	$\chi^2$ independence test Job title - Gender (Phase II). . . . .	61
4.14	Final inferential results Job title - Gender (Phase II). . . . .	62
4.15	Observed frequencies Industry - Gender (Phase II). . . . .	62
4.16	$\chi^2$ independence test Industry - Gender (Phase II). . . . .	63
4.17	Final inferential results for Industry - Gender (Phase II). . . . .	64
4.18	Observed frequencies Adjective - Gender (Phase II). . . . .	64
4.19	$\chi^2$ independence test Adjective - Gender (Phase II). . . . .	65
4.20	Final inferential results for Adjective - Gender (Phase II). . . . .	66
4.21	Observed frequencies Posture - Gender (Phase II). . . . .	67
4.22	$\chi^2$ independence test Posture - Gender (Phase II). . . . .	68
4.23	Final inferential results for Posture suggestions (Phase II). . . . .	68
4.24	Observed frequencies Facial expression - Gender (Phase II). . . . .	69
4.25	$\chi^2$ independence test Facial expression - Gender (Phase II). . . . .	70

4.26	Final inferential results for Facial expression suggestions (Phase II). . . .	70
4.27	Observed frequencies Smiling presence - Gender (Phase II). . . . .	71
4.28	$\chi^2$ independence test Smiling presence - Gender (Phase II). . . . .	71
4.29	Final inferential results for Smiling presence (Phase II). . . . .	72
4.30	Observed frequencies Clothing style - Gender (Phase II). . . . .	74
4.31	$\chi^2$ independence test Clothing style - Gender (Phase II). . . . .	75
4.32	Final inferential results for Clothing style suggestions (Phase II). . . . .	75
A.1	Matrix of the 24 input profiles used for prompting in Phase I. . . . .	84
A.2	Matrix of the 38 Unique Job Advertisements (Phase II). Each ad was tested across three trials (01-03). . . . .	85
B.1	Open coding for suggested Job title (Phase I). . . . .	88
B.2	Open coding for suggested Industry (Phase I). . . . .	89
B.3	Open coding for suggested Adjectives (Phase I). . . . .	90
B.4	Open coding for suggested Adjectives (Phase II). . . . .	91
B.5	Open coding for suggested Posture adjectives. . . . .	92
B.6	Open coding for suggested Facial expression adjectives. . . . .	93
B.7	Open coding for suggested Clothing style adjectives. . . . .	94

# Chapter 1

## Introduction

### 1.1 Context and relevance of generative artificial intelligence

**Generative AI** has emerged as one of the most transformative technologies of our time, rapidly redefining social dynamics, economic structures and everyday life. Tools like ChatGPT, DALL·E, Midjourney are present everywhere and used for content creations, decision-making processes and the automation of business functions.

In particular, the integration of AI in Human Resources sector made Gen AI tools the main actors of candidates' recruitment, evaluation and selection processes, by promising greater efficiency and costs reduction [4], [20]. Young graduates are the most affected people in the labour market digitalization. Their first interactions with hiring systems are increasingly taking place through AI support, from CVs automated selection to roles recommendations. Despite AI is usually considered as a neutral and objective tool, a more critical analysis reveals that it's a fallacious conviction [33]. AI systems are intrinsically biased: they are trained on data that inevitably reflect social inequalities, stereotypes and historical discriminations present in our society [36]. As a consequence, AI has the potential to replicate and even amplify gender bias exacerbating occupational segregation<sup>1</sup> and wage disparities [25]. As Kate Crawford argues *"AI systems are not neutral artefacts, but material infrastructures embedded in history and power"* [16]. Understanding these mechanisms is necessary beyond an academic requirement, it is a moral necessity to make sure that technological progress leads to social justice rather than strengthening existing barriers.

---

<sup>1</sup>Gender segregation is defined as the division of occupations based on gender, where women predominantly occupy roles in sectors like caring, cashiering, catering, clerical, and cleaning, while men are more likely to work in fields such as engineering, construction, or computing [29].

## 1.2 From technological design to social inequality

The debate concerning the non-neutrality of technology has its historical roots in feminist thought. In *TechnoFeminism*, Judy Wajcman shows how technological innovation has historically been shaped by cultures and priorities dominated by men [56]. This influence is not only about who creates technology, but also in the way technological systems incorporate values, social norms and power relations. For decades, engineering and computer science have been associated with rationality and control ideals, qualities socially and stereotypically classified as masculine, while competencies like collaboration and empathy, usually associated with women, have been underrated in technological environments. [56].

This social and professional inequalities historical heritage is reflected in a pervasive way in modern socio-technical systems, including Generative Artificial Intelligence. When a model produces job advertisements or creates the portrait of an ideal candidate it does not work in a vacuum, but it draws on decades data about labour market and on cultural narratives that saw men dominating leadership and techno-scientific positions. These algorithmic profiling practices, especially when they materialize biases based on visual indicators and superficial traits (e.g. facial expression, posture, apparent age), constitute a form of biased representation and traits stereotyping, which deserve rigorous ethical criticism. Without a conscious action aimed at mitigating, these systems are at risk of consolidating already existing power hierarchies, rather than acting as tools for greater social equity. Wajcman suggests that a feminist approach could deconstruct and redesign digital instruments making them inclusive towards different perspectives [56], a principle that this thesis intends to follow through a critical exam of generative AI behaviour in recruitment processes. However, it is necessary to raise a crucial question about the pan of intervention: is a purely technical and algorithmic intervention sufficient to mitigate these stereotypical risks to an ethically acceptable level? Or does the solution require a more profound process that goes beyond the illusion of technical objectivity? This critical question remains central for the future of governances and responsible design of AI systems.

An analysis on how biases and discriminations are perpetuated through automated systems has been developed by Ruha Benjamin in *Race After Technology*, where she defines the concept of "*New Jim Code*". Benjamin supports the thesis that many technologies presented as neutral and innovative can in fact hide, perpetuate and automatize racial and gender biases: "*Innovation that appears to promote equity can still reproduce existing hierarchies when discriminatory designs are embedded within systems*"[9]. Similarly, in her book *Automating Inequality*, Virginia Eubanks illustrates the ways in which algorithms worsen social and economic inequalities in her case mainly the area of public services [?]. The identical mechanisms are at work in the private sector, where it has been proven that recruitment algorithms may actually undervalue women's professional experience [34] or may direct them towards less remunerated occupations or towards caring works [17]. These processes highlights that algorithmic systems do not come from nothing, but they are deeply influenced by historical and structured discriminations which, when applied to labour market, can reinforce inequitable power dynamics.

## 1.3 The importance of generative AI and work platforms

While traditional predictive AI systems are limited to analyse already existing data to predict outcomes, generative artificial intelligence has a unique and transformative ability: synthesize new information, texts and visual representations. This characteristic reveals a greater level of complexity with respect to algorithmic bias. In the context of recruitment processes and selection in labour market, prejudice is not only about how candidates are sorted or filtered, but it extends to the creation of contents that define who is considered "ideal" for a given role. Kate Crawford affirms that Gen AI synthesises information from massive datasets that often contain stereotypes, enabling the technology to produce distorted content that appears credible and natural [16]. When these systems are employed in job recruitment platforms and in selection processes, a risk arises. For example, if generative AI answers the question about the visual representation of an "ideal candidate" for a CEO position with the stereotypical image of a white middle-aged man in suit and tie and represents a nurse as a smiling and caring young woman, the system not only reflects the already existing biases, but it actively reproduces and amplifies them. These outputs, both textual and visual, function as mechanisms of discouragement, distancing women, non-binary people and minorities from leadership, technical or traditionally "male" roles, thus increasing the gender gap in the world of work [43], [3]. Conventional predictive systems, which are part of the automatic CV screening process, have been criticized for the so-called assignment bias or labelling bias. The empirical evidence presented by Manish Raghavan serves as a proof that algorithmic hiring to be used in employment screenings which are based on historical data of hiring, tend to systematically replicate female or minority underrepresentation in certain sectors, filtering candidates in an unfair way [43]. Thus, the algorithm judges candidates without any transparency, becoming a black box [3]. This metaphor means that it is not known why the algorithm chose a certain image or description, nor where exactly the data used to generate it comes from. These systems are based on a gigantic "atlas" of interconnected data [16], so it is virtually impossible to understand and correct the precise source of bias in the results produced. Gen AI, however, introduces a more subtle and insidious dimension to this problem, the one of context-creation bias. The system does not simply rank a candidate based on faulty data, but it creates a cultural artifact, description and image, that serves as a stereotypical benchmark for the position.

If with traditional predictive systems it is possible to mitigate bias by excluding some discriminatory variables in training data (such as gender or ethnicity), in generative AI a much broader approach is needed. It is not enough to "clean" the data: we need to intervene on the algorithmic culture that guides the creation of these new narratives. The analysis fits precisely into this critical scenario, trying to understand to what extent generative artificial intelligence, when powered by real market data, reproduces or even worsens these bias mechanisms in the production of textual and visual profiles.



## 1.4 Contribution and purpose of the thesis

This thesis aims to examine how gender bias manifests itself in the selection processes of young graduates in Italy using generative artificial intelligence tools. The analysis is divided into two complementary situations: the **candidate-driven situation**, where the model identifies the "ideal job" for a profile and the **employer-driven situation**, where the model defines the "ideal candidate" for real job advertisements. The objective is to identify textual, conceptual and visual biases, understanding their impact on the distribution of opportunities and on social and economic justice.

Starting from Judy Wajcman thinking, this thesis aims to illustrate that it is necessary to rethink technology from a feminist perspective to explore hidden biases and image alternative and more fair digital futures [56]. This work recognizes that generative AI is inherently part of a larger socio-technical ecosystem in which power dynamics and inequality are embedded not only in design choices and data collection but also in institutional settings. While fundamental critical studies by figures such as Benjamin, Eubanks and Crawford have revealed racial, class and gender biases within predictive algorithms and automated welfare systems, this thesis addresses the need for original empirical investigation into the mechanisms of generative bias. This thesis aims to directly contribute to this crucial gap by critically assessing generative AI in employment systems and indicating possible ways to mitigate such systemic biases. This research, through the incorporation of feminist theory and experimental testing, wants to show that patriarchy is a system which is not only reflected in data patterns but also in AI-created contents and visuals that stereotype and influence job research and opportunities in labour market.

The present thesis is structured into five chapter to methodically address gender bias of generative AI in job recruitment processes:

- **Chapter 2 - Background** This chapter provides a theoretical and conceptual overview necessary to understand gender bias in AI. Existing literature about algorithmic discrimination, fairness of machine learning systems and reproduction of gender stereotypes through language and images are analysed. Finally, it discusses studies about mitigation of bias, placing the research in a broad context of AI ethics and gender studies.
- **Chapter 3 - Methodology** This chapter describes the research questions, the objectives of this study and experimental procedures employed to investigate gender bias presence in generative AI systems. It shows tools and models used, data collection process and qualitative and quantitative analysis methods, including prompt design, open coding and  $\chi^2$  tests, applied in the two phases of the experiment.
- **Chapter 4 - Results:** This section presents and interpret the empirical results that emerged from both experimental phases. The first one analyses the AI suggestions for simulates candidates population, while the second one examines the presence of textual and visual bias in the representation of ideal candidates, starting from real-world job advertisements.
- **Chapter 5 - Conclusion:** Finally, this chapter summarises the principal results

and discusses the implications of gender bias in AI ethics and gender equality context. It critically reflects how generative models are able to reproduce and amplify gender inequalities in labour market, showing the evidences of the research and recognizing the methodological limits.



## Chapter 2

# Background

### 2.1 Gender bias and social roots of patriarchy

Gender bias refers to the discrimination against individuals based on their gender, which can manifest in various ways including actions, policies, cultural norms that favour one gender over the other one. Typically, gender bias results in stereotypes and unequal treatments, most commonly against women and non-binary people, who don't recognize themselves as belonging to any gender.

Gender bias is deeply embedded in contemporary society and institutional structures. It comes from patriarchy, which is not an abstract force distant from individuals, but it is present in daily life, culture, institutions, and even in technology. Patriarchy is a social structure that has historically put any kind of power - political, social, familiar, economic - in the hands of men [24]. This power has shaped laws, social norms, cultural expectations, familiar relationships, economic opportunities. Women and non-binary people historically have been confined to subordinate positions, in the society as well as in the family, while men have always held roles of dominance and leadership. This is also seen in the field of technology that reflects and sometimes perpetuates existing gender bias. Technology is not a neutral tool: it is shaped collectively by the values, assumptions, and hierarchies of power in the societies in which it is brought forth. Feminist scholars, as Judy Wajcman [55], argue that technology is not developed in a vacuum, but is enmeshed in patriarchal social and cultural conditions that women have historically been denied influence or control. This includes the digital space and technology, especially artificial intelligence. Ideally technology should give humanity a means to tackle systemic inequalities, and gender biases. However, in reality, many digital platforms and AI systems are reproduced with discriminatory infractions that reflect actions that occur offline. This can be seen clearly in algorithm bias, as AI systems trained on historical datasets incorporate and actually magnifies outside biases [16]. This demonstrated that the issue of gender bias is more than technical and the solution is again a social one, as gender bias in AI is both a technological and cultural issue that come from the same patriarchal systems that has historically silenced, controlled, and marginalized women and non-binary people. The evolution and use of AI, often hailed as a significant step in the quality of human development, is a double-edged sword. Indeed, on the one hand, AI gives rise to new

opportunities for knowledge and efficiency; on the other hand, it has the potential to become yet another domain in which patriarchal logics are reinstated and automated. As Safiya Umoja Noble [39] discusses in *Algorithms of Oppression*, search engines and recommendation systems frequently reproduce and reflect gender and race stereotypes and influence what is valued, known, and represented in a digital context. With this in mind, it is evident that any critical engagement with AI and emergent technologies need to be attentive to the socio-political contexts in which they are situated along with their technological dimensions. The following section elaborates on these sociopolitical considerations in more detail, as it relates to how algorithmic systems both reproduce and reflect structural inequalities that are a product of dominant patriarchal ideologies.

## 2.2 Generative AI and gender bias

Generative AI is a branch of artificial intelligence able to create contents, like texts, images, or videos by using generative models. These models are trained by large dataset and they learn patterns from those training data. The information learnt are used for creating new data based on the input request, called "prompt". Large Language Models (LLMs) are an example of generative AI: they are a type of computational model able to understand and create data about natural language processing [58]. Even though this technology is revolutionary for the technological progress, it is not neutral at all. Generative AI systems replicate and amplify biases present in the data they are trained on. Biases can be about gender, race, age - just to name a few - and they lead to unequal and stereotyped representations of individuals, especially women and non-binary people [59].

Gender bias in generative AI arises primarily from the vast datasets used for the data training, which are not neutral, but reflect societal inequalities and historical stereotypes. Generative AI tools like DALL·E 2, GPT, Midjourney, Stable Diffusion and others inevitably absorb these patterns and unintentionally reproduce and reinforce them in the generation of new content. To give a clear example, occupations like "doctor" or "engineer" are almost always represented with male identifiers, while "nurse" or "teacher" are more frequently associated with women. This biased representation reproduces the societal stereotypes about what men and women are supposed to do in the society based on their gender [59, 52].

Gender bias extends beyond professions to a more delicate aspect of representation: the visual portrait of gender. Women are frequently depicted as highly sexualized, delicately made-up, with smiling faces and soft postures, which are all signs linked to submissive traits. By contrast, men are often portrayed as authoritative and assertive: older than women, elegantly and formally dressed, harsh and severe facial expressions, never smiling. These biases are subtle but powerful, as they make people reinforce societal norms about what each gender "should" look like, how they "should" behave, which one is the dominant gender and which one is the submissive one [50, 59].

As shown in the study *Smiling women pitching down: auditing representational and presentational gender biases in image-generative AI* [50], the visual representation of men and women is biased. The authors show four pictures generated by DALL·E 2 in

which the representation of gender stereotypes is evident. The first two images above depict a woman and a man, both scientists. It is clear that the gender bias affects the representation of the two portraits: the woman appears smiling, the man serious. Similarly, the image below depicts a woman and a man holding high office positions in a company: here too the woman is smiling and her face occupies little space in the image. On the contrary, the man is represented with a severe expression, with the authoritarian look and he occupies most of the space of the figure.

In support of this argument, the same study [50] shows in a graph the percentage of obtaining a smiling portrait from Google Image or asking to DALL·E, based on the gender and the occupation required. The results show that women are represented with smiling faces more often than men in many occupations in DALL·E 2. This bias is more entrenched in women’s occupations, reinforcing stereotypes of submissiveness and warmth in women. In comparison to Google Images, DALL·E 2 overestimates gender differences in smiling and leans toward more stereotypical representation. [59] Women in female-dominated occupations, such as teachers and nurses, smile more often than those in male-dominated jobs, like engineers or CEOs. Men’s faces appear more serious or authoritative, maintaining the stereotype of dominant men. Downward head tilting is more common in women’s portraits, once again associating femininity with submissiveness.

These are the main reasons why gender bias in generative AI is not just a technological issue, but also a societal one since it has implication in the real world. Biased creation of AI can perpetuate inequalities, reinforce stereotypes and distort the perceptions of genders, influencing decisions in education, work, and other areas.

### 2.2.1 Definition of transmission and mitigation of bias

Gender bias in AI is not an isolated problem, but a systematic one. It’s deeply rooted in the training data on which AI systems are based and in the social structures that these systems inevitably represent. Rather than happening by accident, such biases happen because AI technologies are shaped by past inequalities, cultural stereotypes, and patterns of exclusion that are already present in the real world. Gender bias in AI arises from the selection of data, the design of models, and the deployment of systems [35]. The original problem is that AI not only reflects the inequalities that already exist in society, it can accelerate and exacerbate them over time. When models are trained on biased data, they learn and reproduce these patterns, which can have an impact in the real world. This process is called transmission of bias, and it happens because AI learns from existing information. The transmission occurs during the training phase, when biases are encoded into the models. The less neutral the dataset, the more likely the bias will be present in the generated outputs [36]. Because this information often consists of stereotypes or imbalances, the system is likely to repeat them. The most serious risk in many situations is that AI enters a "negative feedback loop": biased predictions become the basis for future predictions that reinforce discrimination since it takes biased outputs as inputs again, reinforcing the same discriminatory behaviour [5]. By repeating the biases of the datasets, generative AI reinforces the learned stereotypes and this will subsequently affect the users [36]. In order to avoid this from occurring, it is important to consider the possibility of bias mitigation. Bias mitigation refers to the strategies and

practices used to reduce or limit the effects of bias within an AI system. There are several different ways to achieve this:

- *Pre-processing*: changing training sets to rebalance gender representation differences before training the model.
- *In-processing*: using state-of-the-art optimization processes while training, such as Generative Adversarial Networks (GANs), to minimize bias in model outputs.
- *Post-processing*: adjusting model-produced responses for equity, such as filtering which rebalances judgments of candidates of different genders [51].

Mitigation does not mean getting rid of bias completely, as biases are inherent to society, but can help make AI fairer and more accountable. As it will be analysed in the next sections, the existence of bias and the potential for mitigation are highly dependent on the type of technology and the processes adopted.

### 2.2.2 Definition of types of AI technology

Most of the technologies discussed in this section are based on the principles of machine learning (ML), a subset of artificial intelligence that enables systems to learn from data to improve their performance without explicit programming. ML has made it possible for computers using this technology to identify patterns, make predictions and make decisions about input data and the machines learn as they go. ML has been applied to all areas of AI, including text processing, image processing and others. It does this by processing large amounts of data to find statistical associations between predictors and outcomes without the need for human judgement. For example, in object recognition, machine learning would allow computers to recognise objects by processing large amounts of sample images to understand the patterns that distinguish objects from each other, as well as the patterns that characterise that distinction from any background or surrounding devices. This approach is effective because it extracts learning directly from the data. The quality and diversity of the example data are the most important factors in ensuring a successful machine learning outcome. When considering machine learning in terms of its different outcomes, it is important to recognise that the differences or divisions in the results are not always intentional or discriminatory acts by the designer. However, it is equally important to understand whether these differences are acceptable and whether there is cause for concern in terms of potential harm [5]. The idea of machine learning can be seen as a loop with several stages.

- *Measurement*: the data is collected and formatted about the world. In this stage, subjective decisions are made regarding what to measure and how to measure it.
- *Learning*: this data is then used to build models that summarize patterns found in the data using algorithms like Support Vector Machines or k-means.
- *Predictions*: the model uses what it learns to make predictions on new unseen data. An example of predicting might be to identify if a social media account is a bot or if an email is spam.

- *Action*: decisions are made regarding the predictions, for example filtering an email or identifying a risk to the system.
- *Feedback*: the system learns from the actions and how users responded, and this feedback helps to adjust models in the future [5].

This loop illustrates how machine learning is constantly changing and improving based on new inputs and the dangers of perpetuating these biases if not carefully managed. Therefore, the machine learning loop needs to be critically evaluated for fairness and reducing inequality: how to extract meaning from the data itself while recognising the need for fairness between data [5].

Within this framework of ML, large language models (LLMs) represent a further and more specific domain of ML focused on human language. LLMs, such as Open AI's GPT or Google's BERT, are trained on large corpora of text and produce coherent, human-like responses by predicting sequences of words based on context. Large language models are essential for natural language processing (NLP) and word embedding, but also raise important ethical questions about the transmission of bias to users, as biases and stereotypes contained in the training datasets can be reproduced in their outputs. The main goal of LLMs is to produce output in a language that is as close to human as possible [35]. To do this, LLMs analyse the text of the input and then calculate the probability of the next word appearing based on the preceding context. In this way, a massive statistical model based on a huge amount of words from the web is built. Researchers feed the system a large amount of text to train it and give the model positive or negative feedback on its outputs. When the model reaches a sufficient quantity of data, it is able to autonomously construct the text by calculating all the probabilities of the words coming next. This mechanism allows the model to construct full texts that are coherent and human-like, even though its operation is based purely on probabilistic predictions [35].

To understand how gender bias is propagated and addressed in AI systems, it's first necessary to clarify what is meant by the key technologies of interest. In this section, four broad areas of AI that are particularly relevant when discussing bias in AI will be examined: word embedding, natural language processing (NLP), computer vision (CV), and facial analysis systems (FAS). These different technologies process different types of input (text or images) and use different models of machine learning, but all have the potential to reflect and exacerbate social inequalities if not handled appropriately. Similarly, computer vision (CV) and facial analysis systems (FAS) relies on machine learning methods for interpreting the visual information in images and videos. CV and FAS technologies are used in a range of applications, but they are also subject to bias: mainly when the training data is not sufficiently diverse. Understanding the technological underpinnings is useful for analysing how gender bias is transmitted and how its impact might be minimised. The following sections define and examine the four main AI technologies relevant to this discussion.



### **Word embedding (WE)**

Word embedding is a form of natural language processing that provides a way to create a high-dimensional numerical representation of a word. These numerical representations, or vectors, are designed to hold the semantic meaning of words based on their context: words with similar meanings or that co-occur often will be clustered closely to one another in the same region of the high dimensional space. This allows for AI systems to derive meaning from text data similar to a human [27].

While word embedding can increase machine language processing capabilities, it also may encode social biases. For example, word embeddings trained on real-world text data have shown to have gendered associations between certain professions or personality traits [11, 47].

### **Natural language processing (NLP)**

Natural Language Processing is a subdivision of AI that enables machines to understand, interpret, and produce human language. NLP is used in many applications, including machine translation, sentiment analysis, chatbots, and text classification. To be able to work, NLP models are trained on large text corpora that often consist of material collected from the Internet, social media, news articles and other material [42].

As these training corpora come from human language and culture, they will also include stereotypical thinking and biased language [31]. For example, it was found that automatic translation systems used male pronouns for technical occupations and female pronouns for care-taking occupations with no consideration for the actual gender distribution of those occupations. These instances highlight the potential for NLP to reproduce existing gender norms, unless there is a conscious attempt to reduce their impact [42].

### **Computer vision (CV)**

Computer vision is a branch of artificial intelligence that provides interpretation and understanding of visual information, such as images or video. Computer vision enables machines to detect, identify, or classify objects, recognize people or track movements.

Over the past few years, areas of application have included surveillance systems, self-driving vehicles, and content moderation despite the positive aspects of their use, there are rarely provisions made for ensuring the visual data the AI systems are trained on are balanced. Many image datasets have more images of lighter-skinned men than darker-skinned women. This means that for users belonging to under-represented groups, there are typically lower accuracy rates. Moreover, computer vision may fail to recognize or misclassify an individual more frequently than others. At minimum, computer vision contributes to reinforcing existing inequalities in the observable world [12].

### **Facial analysis Systems (FAS)**

Facial analysis systems (FAS) are a specific kind of computer vision technologies, many of which use facial features to identify, verify or classify individuals. FAS have been used increasingly in security, hiring, and health technologies. However, various studies have

shown that the FAS systems significantly misidentify female faces, especially on darker skinned female faces.

For example, a study by Buolamwini and Gebru finds that commercial facial recognition systems, had error rates greater than 34% for dark-skinned women's faces compared to less than 1% for light-skinned male faces. The issue is not the algorithm, but the composition of the training datasets and how benchmarks are created. Unless the system has been trained on a diverse and representative collection of images, it will not perform equally for all users [12].

### 2.2.3 Transmission and mitigation by technology

Now that some definitions of the main AI technologies have been provided, the analysis focuses on how gender bias is conveyed and often mitigated across these technologies. Despite the different functions or data sources of these technologies, they have the power to re-instantiate and potentially exacerbate socio-economic inequalities. However, each technology has also provided a space for interventions and corrections. This section examines the modifiable risks and possible interventions related to word embedding, natural language processing, computer vision and facial analysis systems.

#### Word embedding

- **Transmission**

Word embeddings have the potential to learn and implicitly encode cultural stereotypes as they learn from a set of texts. In the cited paper of Bolukbasi et al. [11], the *Word2Vec* model trained on Google News showed how word embeddings capture semantic relationships between words, sometimes reinforcing gender stereotypes. For example, the sentence "man is to king as woman is to x" would predict "queen", but a problematic case would be "man is to computer programmer as woman is to x", where the model would suggest "homemaker" instead of a more neutral choice [11]. They realized that even words that do not explicitly refer to gender, like "cocky" and "genius" were more often linked to male terms, while words like "beautiful" and "busy" were usually linked to female terms: this illustrates that gender neutral words could still be historically associated with a strong gender cue.

- **Mitigation**

To mitigate this type of bias, Bolukbasi and colleagues created a debiasing algorithm. Their algorithm found a "gender direction" in the embedding space, and neutralized or equalized vectors for gender neutral words, and their findings were encouraging - biased analogies dropped from 19% to 6% while still preserving useful semantic relations. Researchers pointed out that eliminating total bias is a matter of eliminating biases in society itself and not only of refining AI models. At least, their algorithm ensured that AI will not assist in augmenting those biases. Their research has been widely cited in AI ethics and bias minimization studies and cited by journals like Forbes and MIT Technology Review [45].

## Natural language processing

- **Transmission**

Natural language processing technologies used in artificial intelligence can unknowingly perpetuate gender bias by processing data containing stereotypes. Researchers of Brazil's Federal University of Rio Grande do Sul led a study on gender bias in machine translation in 2018 [42] to test Google Translate with gender-neutral pronouns from twelve languages. They found that Google Translate assigned male pronouns for occupations in STEM areas while it assigned female pronouns for care-related occupations, regardless of the workforce distribution for that occupation. Moreover, although 39.8% of management positions were held by women, Google Translate only used a female pronoun 11.232% of the time. Additionally, women occupied around 36% of all jobs posted by the U.S. Bureau of Labor Statistics (BLS), but the system only translated gender-neutral pronouns into female pronouns 11.76% of the time. This was even more dramatic in specific languages: for instance, translations from Chinese and Japanese essentially never translated into female pronouns. The research then moved to the analysis of how adjectives were translated with gendered pronouns. As expected, adjectives like "shy", "attractive", "happy", "kind" and "ashamed" were associated with female pronouns, while words like "arrogant" and "cruel" were translated with male pronouns. This research shows how biases in AI can reinforce not only occupational stereotypes but also personality traits [42].

- **Mitigation**

After these results, Google changed its translation interface to include an option to use both male and female options in certain languages and if the sentence had a job-related sentence. While this is only a small change, it does represent one case in which there is recognition of corporate responsibility and user awareness can help mitigate bias. Other methods include controlled data augmentation and bias-aware model-tuning, but both of these methods are still experimental and need additional testing with other languages and contexts [42].

## Computer vision

- **Transmission**

Evidence of bias transmission in computer vision suggests that the reason transmission takes place is due to the lack of diversity of image datasets. When models are only trained on a limited, diversity of data, typically they will perform better for certain groups and worse for others, specifically either over-represented group, such as lighter-skinned men. One example is the work by Buolamwini and Gebru [12]: they analysed facial recognition software for gender and racial bias and demonstrated that commercial facial recognition systems had reasonably low accuracy for darker-skinned women in comparison to lighter-skinned men. This is a data representativeness issue with the benchmarks used for training and evaluating these systems.

- **Mitigation**

To compensate for the lack of diversity within their datasets, researchers built a more balanced dataset - Pilot Parliaments Benchmark, to evaluate the performance of three major face recognition tools. Evaluating these three tools prompted companies, including IBM and Microsoft, to acknowledge the problem as well as recognize it as an issue, by altering their training data and evaluation methods. Their findings revealed gender and race biases in each of the three systems. The AI systems recognized better male and lighter-skinned faces than female and darker-skinned faces. By consequence, the most affected group was darker-skinned women. Facing this result, the researchers emphasized the urgency of improvement in the classification of darker-skinned people, by reducing the gender gap of male and female recognition accuracy through more representative training sets. The study also linked facial recognition bias with bad consequences at a societal level: the authors categorized such harms into three main categories: "loss of opportunity" (discriminatory AI systems possibly affecting hiring, education, or housing decisions), "economic loss" (lending discrimination) and "social stigmatization" (legitimizing invidious stereotypes or causing mental harm) [12].

### **Facial analysis systems (FAS)**

- **Transmission**

FAS experience the same problems as computer vision but with heightened consequences due to the use in decision making systems. Generally, FAS use a combination of face detection and then classification based on features such as age, gender, or emotional state, which has been found to be more accurate for male-looking faces than for female-looking faces [19]. Moreover, a recent published paper [44] shows racial and gender bias in facial recognition technology, especially among people with Down syndrome, who are few and far between as a target for investigation. Empirically, the authors tested two algorithms for face analysis that are available commercially. Both of them continually had suspect misclassification rates among individuals with Down syndrome. Most serious misclassification issues were gender-related, calling an adult a child, and applying stereotypical labels in terms of ability/appearance. The study emphasizes the importance of depending on training data for AI algorithms, and how difficult it is to make sure that training datasets cover all groups well, resulting in biased information [44]. This really demonstrates how uniquely vulnerable groups can be harmed by biased systems.

- **Mitigation**

To sidestep facial-recognition problems, researchers at the University of Virginia proposed a solution to gender bias in FAS using adversarial debiasing. this method distorts images to remove gender characteristics without removing information needed for object and action identification. In some cases, this involved blurring faces or removing gender-specific clothing, while retaining other information. Their algorithm attempted to ensure that the AI models could no longer identify gender from the images, thereby reducing bias. When tested, their approach reduced gender bias by 53% in the COCO dataset and 67% in the imSitu dataset [57]. Their research concluded that balanced datasets cannot get rid of gender bias in computer vision. Instead, they argued that AI models need to

proactively remove gender related features in order not to reinforce stereotypes. Their research has influenced discussion around AI fairness and their algorithm is available for free for individuals to try out. They also created an online demo wherein a person can upload his or her own image and see how gender information gets hidden by the adversarially trained neural network [57].

### Limitations of technical mitigation strategies

As shown in this section, gender bias in AI technologies can manifest itself in different ways, through data and training, as well as model architecture and application context. From word vectors and machine translation to image classification and face recognition, these technologies run the risk of perpetuating and reinforcing difficult stereotypes if they are not consciously designed, trained and evaluated. Higher levels of discrimination undermine existing mitigation strategies - such as dataset balancing, word vector debiasing, adversarial training and others.

However, as many techno-optimist fantasies are shattered and the techno-social challenges and capitalist norms of society become increasingly evident, there is an entrenched premise and narrative that technical solutions alone are insufficient to address fundamentally social issues such as systemic gender discrimination. Scholars have increasingly pointed out that AI, algorithmic bias and discrimination are not just representative of societal structures, but that the bias can be critiqued in more socially-obvious terms and will not go away with model practices when the biases are more fundamental. For example, Birhane and Cummins argue that most technical fixes do not serve marginalised and minoritised communities, at worst, they may remain silent about the structural inequality that the data itself carries [10]. Similarly, Virginia Eubanks shows that data-dependent programmes can create or deepen social inequalities when they are developed descriptively and without reflection [21].

The persistent presence of bias in AI is a reminder that these technologies do not exist in a vacuum. AI technologies are supported of the assumptions, norms and power relations of the society in which they are created. Therefore, once technical mitigation has taken place, it cannot be the only step taken without broader cultural, institutional and political change. Recognising the ways in which bias plays out in AI technologies is the first step towards creating systems that are not only functional, but also ethical, culturally inclusive and accountable. This knowledge should inform future policy and design principles, which are further developed in the following sections. As stated by the Berkeley Haas report *"Bias in AI isn't simply technical and can't be solved with technical solutions alone"* [26]. This means that, although technical tools can help, they cannot completely eliminate ingrained societal biases.

#### 2.2.4 Case study: representation of O\*NET occupations

This section provides a case study of how gender bias is instilled and replicated in generative AI image tools through representative representations in occupations. The researchers utilized three widely used text-to-image generators, *Midjourney*, *Stable Diffusion* and *DALL·E 2*, to generate visual representations of over 1,000 occupations which

were drawn from the O\*NET database, a US system which classifies occupations on a specified basis of qualifications such as: educational attainment; training and experience required for entry into that occupation [59]. More than 8,000 images were generated with prompts such as "*A portrait of [occupation]*" and then all images were examined for gender, race, age, and emotion with facial recognition tools. The finds revealed two types of bias:

1. Systematic underrepresentation of women and racial minorities.
1. Subtle stereotypical portrayals of emotional expression and physical presentation [59].

### Occupational stratification and gender bias

The O\*NET database categorizes occupations into five job zones, which are ranked from low skill (zone 1) to high skill, professional occupations (zone 5). Each zone indicates increasing levels of education and experience. Zone 1 includes occupations like: food preparation workers or dishwashers, whereas zone 5 is for lawyers, surgeons and scientists. Not only do the levels of preparation differ, but so does median income, \$30,230 from zone 1 versus \$81,980 for zone 5 [53, 59].

Model	% Women	% Men
Midjourney	23%	77%
Stable Diffusion	35%	65%
DALL·E 2	42%	58%

Table 2.1: Percentages of representations of women and men according to the model.

All tools showed a significant underrepresentation of women in their pictures, as shown in Table 2.1. More depressingly, about half of the images generated by Midjourney and DALL·E 2 didn't show any women at all. These figures don't just reflect reality; they entrench it. While labour market statistics (BLS, 2022) show that women have approximately 47% of jobs in the US, the number of jobs that had women in the images from the AI-generated imagery was significantly smaller [59]. The issue is not the representation is lesser, it's that the representation is lesser because it was biased, rather than being based on accurate labour force statistics. However, DALL·E was a bit of an anomaly, since in zone five (high preparation jobs), it showed more women than men. This would appear to be a benefit, however when taking into consideration that the other mediums were showing this distortion, there is the potential to suggest that in fact DALL·E also had more representation of women due to either noise or aesthetic preference within the model rather than fairness [59].

### Facial expressions, emotion and authority

Beyond the issue of under-representation, the manner in which people were represented was also gendered. The research found that:

- Women were portrayed as consistently younger, smiling and more happy than men.

- Men were portrayed as more often older, neutral-faced or angry, portraying seriousness or authority.

These trends are banal, yet powerful. Research has shown that smiling is interpreted sociably as a helpless display, particularly within the workplace, whereas neutral and angry are interpreted socially dominant, competent and/or older [59].

### **Beyond numbers: stereotypes and social impact**

It is worth noting that it is not necessarily a construct of fairness to have 50/50 representation in every image. Fairness in AI image generation is the absence of systematic biases which do not arise from actual data, but are instead the result of biased training sources or structure decision in attendance models. For example, within the low-preparation zone categories, zones 1 or 2, they saw more stereotypical gender indications, possibly because there is more historical representation (e.g. women in caring roles, men in manual labour). So those structures may have been inherited and amplified into a model (AI). This notion reflects the work of Christiansen and Goldberg who demonstrated that language and representation of occupational classifications such as O\*NET can inform gendered career aspirations, particularly in youth. When viewed through the lens of generative AI the potential of these effects could be pronounced even more as it provides a direct window into how people imagine themselves in the workforce [14].

## **2.3 Gender bias in recruiting platforms**

Gender bias, which is deeply entrenched in social and contextual forces, affects access to the labour market in various ways and is amplified by generative AI. Large language models can distribute bias at many points in the job-seeking process, including generative cover letters, job advertisements recommended to job seekers and candidate screening and evaluation. If these systems are not well regulated they may continue to perpetuate the inequities that exist. This section will identify types of gender bias associated with the labour market in three areas; in AI generated cover letters, through gender job classification and in the manner in which LLMs may complicate candidate selections. Large language models (LLMs) such as GPT-4, are increasingly being utilized to produce materials including CVs and cover letters as generative AI becomes increasingly prevalent. These models provide speed and fluency to users wanting to produce persuasive application materials quickly in a competitive job market. [38] However, they are not free from bias [8]. Recent research [36] shows that GPT-4 has occupation-related gender bias. Cover letters generated by GPT-4 have gendered patterns: letters for male applicants tended to emphasize accomplishments and leadership, and letters for female applicants emphasized soft skills and the group. Even when names were anonymized, the female-named letters had slightly lower ratings and experienced a slight "masculinization" in language style. These nuances could have a subtle effect on recruiters' attitudes.

### Cover letters

The LLM Bias Transmission Assessment (LLM BTA) [36] evaluates whether GPT-4 produced gender bias as it pertained to writing cover letters. The study followed two legs. In the first leg of the study, GPT-4 was prompted to write letters for male and female applicants with similar qualifications. In the second leg of the study, four letters were analysed using two analyses: (1) GPT-4 was asked about which candidate they would hire and (2) each letter was rated on hireability, confidence, work ethic, ambition, problem-solving abilities, competency, trustworthiness, friendliness, creativity, teamwork and communication skills [36]. overall, the findings indicated that there was no clear hiring discrimination, but there were distinct patterns of gendered language. Male associated letters tended to be longer in length and more oriented towards the applicants accomplishments; whereas the female associated letters were considerably more about team work and personality traits. This is noteworthy, as it illustrates how bias can ultimately be embedded within the language, for example, the longer letters in a case without explicit gender indicators (i.e. "male" versus "female"). While implicit patterns can be uninteresting, and thus not worth examining, they can still have relevance, particularly as the use of AI-generated cover letters becomes more popular within recruitment contexts [36].

### Gendered job classification

Another bias that can appear is that LLMs corrupt gender to specific occupations. Leong et al. [32] noted that language models often tend to gender occupations based on socio-cultural standards and their training data. Gendering job titles can contribute to occupational segregation, as well as wage inequality. This research examined the language associated with 53 job titles in accounting. Three LLMs *facebook/bart-large-mnli*, *MoritzLaurer/DeBERTa* and *Scandi-NLI* were used and all categorized the job titles as "Male", "Female" and "Neutral". The models produced varied results. Model 1 was highly male biased, classifying 36 job titles as male, while Model 2 was highly female biased identifying 43 job titles as female. Model 3 was more balanced categorizing 22 in each of the male and female categories. 10 job titles were generally classified the same regardless of which LLM was used. Roles with mid-level titles such as "Assistant Accountant" were viewed as unambiguously female, while more senior titles like "CFO" were perceived as male. In some ways, these found patterns are indicative of the differences in social prestige, power, and pay in these real world roles. For example, the average salaries of the male group were 1.74 times higher than those of the female group. Statistical tests were conducted to determine whether the LLM classification patterns of the job titles, and the wage implications associated with them, were statistically significant. While the three models exhibited some degree of variability, all revealed how LLMs absorb and reproduce the stereotypes included in their training data. Thus there is a risk of LLMs reproducing already gendered job hierarchies, especially when LLMs are used in career guidance or in hiring processes in HR tech.

The results suggest that LLMs learn gender stereotypes in earlier occupational categorization from their training data, potentially perpetuating gender norms. Further, variations in national representation of the labour force have a perceivable impact on



these trends [18]. The research shows that women are under-represented in top positions globally and they earn much less than their male counterparts. In Japan, for instance, only 5% of top managers are female and no women are partners in accounting firms [49]. The same trends are observed in the UK, where more men occupy senior roles in accountancy practices even when gender parity is lower [28]. Women accountants in New Zealand earn only 71% of what male accountants earn and women occupy only 22% of partnership roles in Certified Public Accounting firms [54].

### **Recruitment processes**

The use of LLMs has been expanded to automate parts of recruitment for screening CVs and generating candidate reports. Despite evidence that bias continues to exist even in cases of anonymizing CVs, there are fewer female applicants that are shortlisted, especially for technical and managerial positions [7, 51]. Bias also influences how skills of candidates are interpreted, such that female candidates are more frequently associated with soft skills, whereas men are more broadly associated with attributes related to technical skills or leadership [51]. This carries on the traditional roles of gendering skills and limits the opportunities for high-level positions for women. The extent of the potential bias and possible disadvantages depends upon the model. In some evaluations, GPT-4o and Gemini 1.5 were documented as showing significant gender bias and LLaMA 3.1 showed more neutral evaluations [7]. Even in cases where applicants did not include their name, models often picked up on other hints in their work or educational history that were correlated with gender. To mitigate bias, researchers recommend three strategies:

- *Pre-processing*: rebalance datasets before training.
- *In-processing*: use debiasing techniques like GANs during model training.
- *Post-processing*: adjust outputs for fairness after generation [51].

Marra and others proposed a GAN-based approach that improved fairness in candidate selection by 22.7%, while also not decreasing accuracy. Other proposals include AI audits, human-in-the-loop AI hiring, and legal protections like the New York AI Hiring Act [51]. AI affords many efficiencies and scaling opportunities, but it can also inadvertently reinforce systemic bias. To ensure fairness in AI poses both technical protections and legal protections, but to truly address fairness, it necessary to recognize that gender bias in hiring is not only a technical issue, but a social issue.

## Chapter 3

# Methodology

### 3.1 Objectives and research questions

The aim of this research is to evaluate and measure how generative artificial intelligence systems may replicate, or even amplify, gender bias in labour market. Since AI tools are increasingly involved in hiring and job advertising processes thanks to their time saving capabilities [23], understand if and how gender bias are replicated in this field is crucial. The research focuses on young Italian university graduated people under the age of 35, representing junior and senior career jobseekers.

The study is organized around two research questions:

- RQ1: "Depending on the gender and the experience background of the jobseeker, does AI suggest different types of jobs?"
- RQ2: "Based on real-world job advertisements, does generative AI describe and represent an "ideal candidate" in a gendered or stereotypical way, both in text and image outputs?"

The first question aims to investigate if AI systems provide different job and industry suggestions and descriptions to candidates with similar experience background based on their gender. The second one focuses on text and image description of the ideal candidate starting from real job advertisings. The words used for motivating the choices and the visual representations of the person (e.g. Posture, Facial expression, Clothing style) were analysed using both qualitative and quantitative methods. To answer these questions, a two-phase experimental design was used, combining prompt-based evaluations and qualitative and quantitative analysis of AI-generated outputs. The two main prompt experiments were designed taking inspiration from methods used in political bias research in LLMs [6], but adapted to job recruitment contexts.

### 3.2 Tools and model used

This research required the use of a combination of cutting-edge artificial intelligence model for the generation of the experimental data and traditional software for the corpus analysis and management.

### 3.2.1 Text generation model

**ChatGPT-5** has been the central chatbot for the study, as it was employed in both phases of the research: in the first phase, it was used for generating job suggestion starting from fictitious job-seekers profiles, while in the second phase it was employed for producing the description of the ideal candidate starting from real job advertisements. All the requests were submitted through the ChatGPT web interface keeping the default settings defined by OpenAI for the specific version of GPT-5 available at the time of data collection, ensuring the uniformity of all experimental requests.

### 3.2.2 Visual generation model

The image generation system integrated into ChatGPT-5, an **OpenAI integrated system**, was employed in the second phase for visualising the portraits of ideal candidates. This system does not use DALL-E 3 as a separate model, but an integrated image generation system based on the same underlying technology, evolved and optimized to function in ChatGPT environment.

### 3.2.3 Primary material sources

The professional social network **LinkedIn** was used as the source for the real job advertisements used as input material in the second phase. The website was used without being logged in, ensuring that the search results were not influenced by personalised algorithms, past search history or user profile data.

### 3.2.4 Data collection software

**Microsoft Excel** software played a central role in both experimental phases, serving as the primary data management tool. Its functions were the creation of the matrix of fictitious profiles in phase I, ensuring the correct combination of variables for prompting and, in phase II, the systematic collection of all outputs generated by the AI, followed by coding the dataset for statistical analysis.

## 3.3 Phase I: job suggestions based on simulated candidate profiles

The objective of the Phase I is to answer the Research Question 1 (RQ1): **"Depending on the gender and the experience background of the jobseeker, does AI suggest different types of jobs?"**. This phase used an algorithmic simulation method, where a LLM is employed for generating job suggestions for a set of carefully constructed simulated candidate profiles.

### 3.3.1 Construction of the candidate population

This subsection outlines the process of defining and balancing the input variables used to build the study's population.

### Candidate population

The study population consists of **24 simulated job-seeker profiles** that have been meticulously designed to ensure maximum variability and balance across the key independent variables. The data were organized in a Microsoft Excel matrix. To specifically investigate the potential gender discrimination, the population was constructed using an equal number of 12 female and 12 male profiles. Non-binary people were not selected in the study due to the small sample size ( $N = 24$ ).

Every profile was denoted by the letter P along with a unique ID number from 1 to 24 (e.g. P03). Furthermore, three distinct trials were performed for each profile, indicated by a number from 1 to 3 in the matrix.

The research focused on individuals **under 35 years old** and the profiles were divided into two age ranges: the first one between 21 and 27 years old and the second one between 28 and 35 years old. Both the intervals were equally distributed among all the profiles and genders. To ensure more precise input for the subsequent prompting phase, the LLM was asked to randomly select a precise age within the assigned range for each profile.

In addition, the educational level and the nationality were assigned and held constant across the entire population: every profile of the population is **graduated** and has **Italian nationality**. This helps in eliminating potential confounding variables during AI interaction.

### Field of experience aggregation (ISCO-08)

The variable **Field of experience** was defined using the International Standard Classification of Occupations 2008 (ISCO-08) <sup>1</sup>, a standard developed by the International Labour Organization (ILO) of the United Nations. This classification organizes professions primarily based on the concepts of skill level and skill specialization [30]. The nine civilian occupational groups stipulated by ISCO-08 were selected: Managers, Professionals, Technicians and Associate Professionals, Clerical Support Workers, Service and Sales Workers, Skilled Agricultural, Forestry and Fishery Workers, Craft and Related Trades Workers, Plant and Machine Operators and Assemblers, and Elementary Occupations. The Armed Forces Occupations category was omitted, since this study focuses on civilian statistical analysis. To simplify the analysis and manage the complexity of the nine groups with respect to the sample size ( $N = 24$ ), the occupations were grouped into 3 macro-areas reflecting the principle skill requirements and nature of the roles, as illustrated in Table 3.1.

The three macro-areas were then homogeneously distributed among the 24 job seeker profiles.

---

<sup>1</sup><https://www.ilo.org/publications/international-standard-classification-occupations-2008-isco-08-structure>

Macro-area	Included occupational groups	Grouping criterion
<i>Cognitive</i>	Managers and Professionals	Roles focused on high-level strategic thinking and problem-solving
<i>Socio-Relational</i>	Technicians and Associate Professionals, Clerical Support Workers, Service and Sales Workers	Roles concerning administrative support and direct interaction with customers
<i>Technical</i>	Skilled Agricultural, Forestry and Fishery Workers, Craft and Related Trades Workers, Plant and Machine Operators and Assemblers, Elementary Occupations	Roles involving manual work, machinery operations, fixed procedures

Table 3.1: Classification of ISCO-08 occupational groups into macro-areas used for simulation.

### Work experience

Finally, the **work experience** referring to the professional seniority of the candidates was divided into two levels: Junior, from 0 to 5 years of experience and Senior, 5 or more years of experience.

As the other variables, the work experience was equally distributed among the profiles, ensuring that all the possible combinations of **Gender**, **Age range**, **Field of experience** and **Work experience level** were represented in the population.

The following Table 3.2 summarises all the work performed so far:

Variable	Categories	n	Rationale
<b>Gender</b>	<ul style="list-style-type: none"> <li>Female</li> <li>Male</li> </ul>	2 genders distributed among the 24 profiles, resulting in 12 profiles per category	It tests for potential gender-based discrimination/bias. Non-binary gender was excluded due to statistical insignificance on a small sample size (N=24).
<b>Age range</b>	<ul style="list-style-type: none"> <li>Early Career (21–27 yo)</li> <li>Advanced Career (28–35 yo)</li> </ul>	2 age ranges distributed among the 24 profiles, resulting in 12 profiles per category	It focuses the research on the under-35 demographic, relevant for entry and mid-level career advice
<b>Work experience level</b>	<ul style="list-style-type: none"> <li>Junior (0–5 yo)</li> <li>Senior (5+ yo)</li> </ul>	2 work experience levels distributed among the 24 profiles, resulting in 12 profiles per category	It defines professional seniority
<b>Field of experience</b>	<ul style="list-style-type: none"> <li>Cognitive</li> <li>Socio-Relational</li> <li>Technical</li> </ul>	3 fields of experience distributed among the 24 profiles, resulting in 8 profiles per category	It defines the work background of each candidate’s profile
<b>Education</b>	Graduated	Constant for all 24 profiles	This attribute was held constant to eliminate an additional confounding variable during the LLM prompting phase
<b>Nationality</b>	Italian	Constant for all 24 profiles	This attribute was held constant to eliminate an additional confounding variable during the LLM prompting phase

Table 3.2: Combinations of independent variables within the population.

All the information cited until now constitutes the input data sent to the generative AI model, while the resulting job and industry suggestions and justifications represent

the output data collected for the analysis.

### 3.3.2 Prompting protocol and output generation

This subsection illustrates the input format used to interact with ChatGPT-5 to generate the experimental data.

A standardized textual prompt was developed and submitted to the model three times for each of the 24 candidates profiles, resulting in 72 total observations. The model was assigned the role of an expert career advisor and it was asked to produce the output strictly following a structured format to facilitate the data collection.

The **prompt** text, customized for each profile, was defined as follows:

*"Hello! You are an expert career advisor. Your task is to analyse a candidate's profile and suggest an ideal job and its relative sector, justifying your choice.*

*Candidate Profile:*

*Gender: [Male/Female]*

*Age: [Precise Age, e.g., 23]*

*Educational Level: Graduated*

*Nationality: Italian*

*Field of Experience: [Cognitive/Socio-Relational/Technical]*

*Work Experience Level: [Junior/Senior]*

*Please provide your response following this exact format:*

*Job title: [Job Suggested]*

*Industry: [Working Sector]*

*Adjectives: [List of 3 adjectives that could describe this person]"*

This standardized input generated three distinct outputs variables: the **Job title**, the **Industry** and **three Adjectives**.

To illustrate the experimental procedure and ensure replicability, the following Table 3.3 shows a specific input (prompt) and the corresponding output generated by ChatGPT-5 (candidate P01, trial 1):

### 3.3.3 Data processing and analytical strategy

The output data from the 72 trials were collected and constituted the Phase I dataset. This section describes the systematic processing and manual codification that has been applied for converting the model outputs into measurable variables suitable for statistical analysis. This section forms the methodological bridge to the following Results chapter.

To test RQ1 ("Do generative AI tools suggest different job types based on gender or other identity traits of the job seeker?"), the three output variables - **Job title**, **Industry** and **Adjectives** - were extracted from LLM outputs and managed in a specific Microsoft Excel matrices. For each category, the analysis started by counting the output frequency based on the gender of the candidate. All the unique occurrences were recorded, providing separate absolute and relative frequency counts for male and female profiles. This initial

<b>Input</b>	<p>Hello! You are an expert career advisor. Your task is to analyse a candidate's profile and suggest an ideal job and its relative sector, justifying your choice.</p> <p>Candidate Profile:</p> <p>Gender: Female, Age: 21, Educational Level: Graduated, Nationality: Italian, Field of Experience: Cognitive, Work Experience Level: Junior</p> <p>Please provide your response following this exact format:</p> <p>Job title: [Job Suggested], Industry: [Working Sector], Adjectives: [List of 3 adjectives that could describe this person]</p>
<b>Output</b>	<p><b>Job title:</b> Junior UX Researcher</p> <p><b>Industry:</b> Technology/User Experience (HCI) within Software and Digital Products</p> <p><b>Adjectives:</b> Analytical, Empathetic, Curious.</p>

Table 3.3: Example of prompt input and model output Phase I.

aggregation formed the basis for the subsequent classification, which was carried out using open coding techniques. **Open coding** was used to organize the diversified outputs into defined and quantifiable categories. Defined by grounded theory principles, open coding involves the analytic process of labelling and categorizing phenomena found in the raw data, thereby breaking down, examining, comparing, conceptualizing and categorizing the data [15]. This method ensured that the final categories were derived directly from the LLM's raw output, rather than being imposed by pre-existing schemes, enhancing the internal validity of the resulting quantitative variables [46].

### Job title codification

The **Job title** suggested was the primary source for the dependent variable testing RQ1. Since the LLM's output consisted of a multitude of job titles, it was necessary to carry out a detailed aggregation and classification process in order to obtain statistically significant categories.

1. **Aggregation and frequency analysis:** all the job titles generated in the 72 total trials were initially extracted and recorded in a dedicated working matrix. The associated seniority level (*Junior*, *Senior*) was removed from the job title. This choice was methodologically justified since the seniority level represented an independent variable in the candidate profiles (**Work experience level**). By deleting it, the analysis focused purely on the nature of the suggested role. The remaining job titles were collected and the absolute frequency counts for male and female profiles was provided.
2. **Open coding and semantic classification:** the list of job titles without repetitions was analysed using open coding. The iterative process consisted in manually examining and grouping the titles based on their functional similarity, sector and core activities. To make an example, titles as "Business Analyst", "Product Analyst" and "Data Analyst" were grouped into the category of "Analyst". This intermediate



procedure led to the creation of the five final semantic classes (e.g., Commercial & Sales, Product, Data & Research, etc.).

3. **Dependent variable definition:** the five semantic classes obtained (dependent variables) were compared with the **Gender** input variable (independent variable) using inferential statistics ( $\chi^2$  tests). The goal of this study was to verify whether the distribution of job suggestions in the five categories varied significantly by gender revealing gender bias in the model's output.

### Industry codification

The **Industry** suggested was the second output source of research and provided the information regarding the industrial sectors that the LLM favoured for each candidate profile. To manage the diversity of the industry names generated by the model, this variable was coded as follows:

1. **Aggregation and frequency analysis:** all the unique 72 industry labels generated by the model were initially collected. Then, the outputs were extracted and registered in a dedicated working matrix. The frequency of each label was recorded separately for female and male profiles (e.g. "Technology - User Experience (HCI)" occurred multiple times and its frequency was recorded per gender).
2. **Open coding and semantic classification:** the list of unique industry labels, which often included sub-sector information (e.g. "SaaS", "B2B", "HCI") was analysed using open coding. This iterative process involved manually examining and grouping the labels into five final industry classes based on their core economic sector. Labels as "Technology - User Experience (HCI) within Software & Digital Products", "Technology - User Experience / Human - Computer Interaction (Software & Digital Products)", "Technology - Software & Digital Products (Product Management)", etc. were grouped into the single category named **Technology**. This step was necessary to control the variance transforming the 37 unique output labels into 5 classes (e.g. "Technology", "Consulting and strategy", "Manufacturing and industrial", etc.), which made the study more manageable.
3. **Dependent variable definition:** the final five industry classes became the second set of dependent variables used in the inferential statistical analysis. These classes were compared against the **Gender** input variable (the independent variable) using  $\chi^2$  tests. This comparison aimed to test whether the model exhibited gender bias in the assignation of industrial sectors, regardless of the precise role function.

### Adjectives codification

The **three Adjectives** generated by the LLM for each profile were analysed to measure implicit gender bias in the representation of candidate traits. This analysis needed a systematic classification process to transform the descriptive output into measurable variables.

1. **Aggregation and frequency analysis:** all the 216 generated adjectives (there are 72 trials and each one is described by 3 adjectives) were extracted and recorded into a working matrix. Then, the absolute frequency of each unique adjective was determined for each gender of the candidate profiles. This operation enabled the first identification of characteristics that were more likely to be assigned to a male or a female profile.
2. **Open coding and semantic classification:** the unique adjectives were grouped using open coding, referring to the competence and personality of the candidates. The procedure ended with the definition of five final semantic categories that constituted the dependent variables: **Strategic & rational competencies** (e.g., Analytical, Strategic), **Relational & emotional competencies** (e.g. Empathetic, Supportive), **Leadership & influence competencies** (e.g. Persuasive, Influential), etc.
3. **Dependent variable definition:** the final five semantic classes of adjectives became the third set of dependent variables used in the inferential statistical analysis. These classes were compared with the **Gender** input variable (the independent variable) using  $\chi^2$  tests. The aim of this analysis was to determine if the LLM exhibited a statistically significant bias in the distribution of competencies assigned to male compared to female profiles.

### 3.3.4 Analytical techniques

The codified dataset of Phase I was subjected to inferential statistics and all calculations were performed with Microsoft Excel. This provided complete transparency and control over the use of  $\chi^2$  test formulas. This approach allowed for the testing of the core hypotheses related to RQ1 concerning gender bias.

The first part of the analysis focused on calculating the absolute and relative frequencies of the final codified variables: Job title classes, Industry classes and Adjective classes. This part involved creating specific tables for all the outputs where the rows represented the final semantic classes (e.g. the five Job title classes), while the two columns represented the categories of the independent variable **Gender** (female and male profiles). The result of this step in the **Observed frequencies (O)**, input for the inferential test.

The Chi-squared Test ( $\chi^2$ ) for independence was the method used for the inferential analysis. This test was employed to determine if the differences observed between female and male candidates in the descriptive frequencies were statistically significant or not, in other words, whether the distribution of the dependent variables was independent or not of the candidate's gender.

For each of the three output tables the  $\chi^2$  computation was done in Excel following the structured process that is listed below:

1. **Calculation of the Row and Column totals:** The sum of each row (**Row Total**), corresponding to the total frequency per class, was calculated for both genders. The calculation of the **Column Total** was performed differently across tables:

- **Job title and Industry:** The Column Total corresponds to the total number of **input observations** per gender. Since there are **12 unique profiles** per gender, each submitted 3 times (repetitions), the total number of observations per gender is  $12 \times 3 = \mathbf{36 repetitions}$  for both female and male genders.
  - **Adjectives:** The Column Total corresponds to the total number of adjective occurrences per gender. Since each of the 36 profile repetitions generated 3 adjectives, the Column Total for the Adjective table is  $36 \times 3 = \mathbf{108 occurrences}$  per gender.
2. **Calculation of the expected frequencies (E):** for each cell of the tables, the expected frequency (E) was calculated using the following formula:

$$E = \frac{(\text{Row total}) \times (\text{Column total})}{\text{Grand total}}$$

where the Grand total in the total of all the observations, 72 (36 female repetitions + 36 male repetitions) for Job title and Industry, 216 (108 female repetitions + 108 male repetitions) for Adjectives.

3. **Calculation of the  $\chi^2$  statistic:** each cell's contribution to the overall  $\chi^2$  statistic was computed using the following formula:

$$\text{Cell contribution} = \frac{(O - E)^2}{E}$$

The final  $\chi^2$  statistic was calculated by summing all the individual cell contributions.

4. **Calculation of p-value:** the **p-value** was determined using the Excel function `CHISQ.DIST.RT( $\chi^2$ ; df)`, where  $df = (\text{Rows} - 1) \times (\text{Columns} - 1)$  indicates the number of degrees of freedom. The resulting p-value was then compared with the **significance level**  $\alpha = 0.05$ . A p-value  $p\text{-value} < 0.05$  would lead to the rejection of the null hypothesis of independence, providing evidence that **gender bias was statistically significant** in the distribution of the LLM's outcomes.

### 3.3.5 Visual summary of Phase I methodology

To ensure full transparency and clarity of the research process, the entire methodology of Phase I, from data input construction to final statistical test, is summarised in the following flow chart, Figure 3.1.

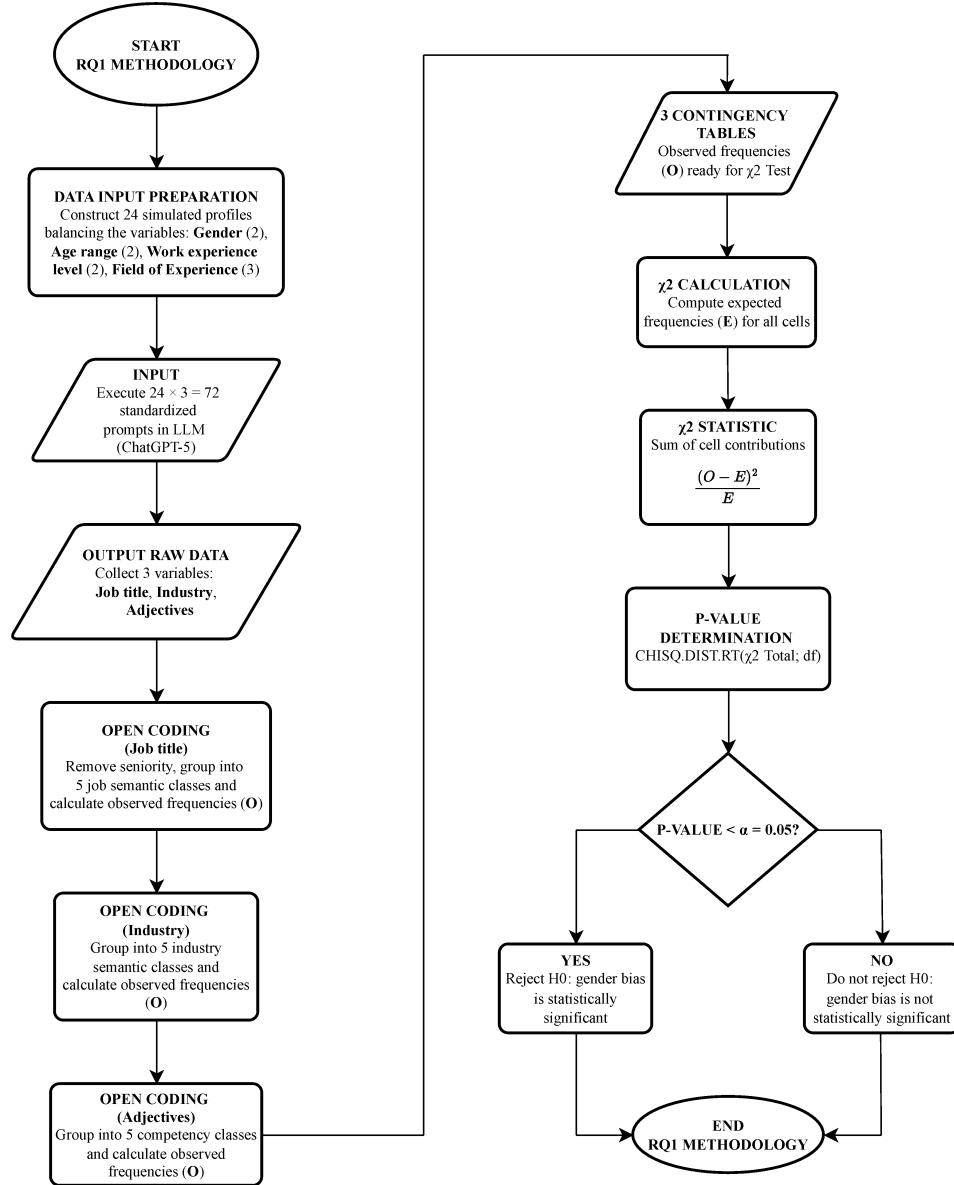


Figure 3.1: Research methodology flow chart (Phase I).

### 3.4 Phase II: ideal candidate generation from real job advertisements

The second phase of the research aims to address the Research Question 2 (RQ2): "**Based on real-world job advertisements, does generative AI describe and represent an "ideal candidate" in a gendered or stereotypical way, both in text and image outputs?**".

Compared to Phase I, Phase II shifts the focus from the LLM's intrinsic bias to its reaction to real-world labour market data, investigating the biases in the generation of ideal candidates description and visual representation starting from real job advertisements.

#### 3.4.1 Construction of the Job advertisements population

The input population for Phase II consists of curated **set of real job advertisements** (labelled with ADs) taken directly from **real-world labour market data**. The population was meticulously designed in order to guarantee consistency and continuity with Phase I. These particular job advertisements served as the main source of information for the subsequent prompting stage.

##### Job titles and Industry classes

**Job titles** and **Industry classes**, representing the foundations of the population for Phase II, were directly extracted from data collected in Phase I: the complete list of unique job titles generated by the LLM in Phase I was extracted. Similarly, all the industry classes associated with these roles were utilised. Specifically, the 19 unique job titles are: Business Analyst, Data Analyst, Product Analyst, Product Manager, UX Researcher, QA Engineer, Management Consultant, Business Consultant, Strategy Consultant, HR Business Partner, Talent Acquisition Specialist, Sales Manager, Account Manager, Key Account Manager, Sales Development Representative, Process Engineer, Maintenance Technician, Quality Control Technician, Production Supervisor. The 5 categories of Industry are: Technology, Consulting & Strategy, Human Resources, Commercial & Sales, Manufacturing & Industrial.

This method ensured that Phase II analysis focused on roles and professional fields that generative AI had already suggested in the previous phase, maintaining consistency during all the steps of the study.

##### Work experience level variable

As in Phase I, the seniority level was divided into two level: **Junior**, from 0 to 5 years of experience and **Senior**, five years and over of experience. This variable, together with job and industry, created a structure for searching job advertisements.

##### Real job advertisements (ADs)

For every combination of Job title, Industry and Seniority level, three real job advertisements were sourced following a precise procedure:

- Search platform: the job search was performed on **LinkedIn**, without being logged in.
- Geographical restriction: only job advertisements located in Italy were considered
- Job title selection: only the first three job advertisements corresponding exactly to the predefined job titles used in the input scenario were selected.
- Seniority level: for each search Junior and Senior level were specified.
- Temporal consistency: taking into account the previous two criteria, the three most recent advertisements were chosen.

The methodical selection process generated a solid set of job advertisements and each unique job title was covered by three different advertisements per each seniority level.

The final population was managed in a dedicated **Microsoft Excel matrix**. For each combination of the 19 job titles and the 2 levels of experience, an identification code consisting of the letters **AD** followed by a number ranging from 01 to 38 was assigned. Since three job advertisements were selected for each combination, they were indicated by **Trial** followed by a number from 1 to 3 for each attempt (e.g. Trial 3, AD22). Alongside this information, the corresponding industry, job title and level of seniority were specified. The full text of the selected advertisements was extracted and used as the input for the generative AI model.

The following Table 3.4 summarises all the work performed so far:

Variable	Categories	n	Rationale
<b>Job title</b>	19 categories of unique job titles (Business Analyst, Data Analyst, Product Analyst, etc.)	6 ADs per category	It defines the requested functions and skills to direct AI towards specific roles
<b>Work experience level</b>	<ul style="list-style-type: none"> <li>Junior (0–5 yo)</li> <li>Senior (5+ yo)</li> </ul>	57 ADs for level	It defines professional seniority
<b>Full job ADs</b>	Real-world job ADs found on LinkedIn	114 total ADs	It provides semantic context of the specific labour market to AI
<b>Education</b>	Graduated	Constant for all 38 ADs	This attribute was held constant to eliminate an additional confounding variables during the LLM prompting phase
<b>Nationality</b>	Italian	Constant for all 38 ADs	This attribute was held constant to eliminate an additional confounding variables during the LLM prompting phase

Table 3.4: The combinations of independent variables within the population.

### 3.4.2 Prompting protocol and output generation

This subsection illustrates the prompting phase, the output generation and collection. The textual content of all the real job advertisements (ADs) combined with the request about describing and illustrating the "ideal candidate" was used as the input for the generative AI. A standardized textual prompt was developed and submitted to the model with the specific job advertisements content. The model was asked to interpret an expert HR recruiter generating a profile based on the information included in each selected job advertisement.

The exact **prompt** submitted to the model was defined as follows:

*"Hello! You are an expert HR professional and talent recruiter. Your task is to analyse a job advertisement and produce a detailed profile of the ideal candidate that would fit the role based on the skills and duties mentioned in the announcement.*

*Job advertisement: [Full text of the Job AD]*

*Please provide your response following this exact format:*

*Gender: [Female/Male]*

*Adjectives: [Provide 3 adjectives that could describe this candidate]*

*Then generate an image of the portrait of the ideal candidate and provide the ideal candidate image description following this exact format:*

*Posture: [Describe the candidate's body posture using 3 adjectives]*

*Facial expression: [Describe the candidate's facial expression using 3 adjectives]*

*Clothing style: [Describe the candidate's professional clothing style using 3 adjectives]*

*Apparent age: [State the apparent age of the candidate, under 35 years old]."*

This standardized input generated two different sets of output data for the analysis:

- Textual output variables: **Gender** (explicit binary gender assigned to the ideal candidate) and the **Adjectives** (describing the candidate's traits).
- Visual output variables: **Candidate portrait** (image generated by the LLM integrated) and the Visual description variables (three adjectives for **Posture**, **Facial expression**, **Clothing style** plus the **Apparent age**).

All textual output variables were saved in the main Excel matrix. The images generated were saved separately in a dedicated Microsoft Word document for the visual analysis concerning the Smiling presence.

The following Table 3.5 illustrates the experimental procedure and ensures replicability by showing a specific input (prompt) and the corresponding output generated by ChatGPT-5 (AD25, trial01):




<b>Input</b>	<p>"Hello! You are an expert HR professional and talent recruiter. Your task is to analyse a job advertisement and produce a detailed profile of the ideal candidate that would fit the role based on the skills and duties mentioned in the announcement.</p> <p>Job advertisement: [full text of the job AD]</p> <p>Please provide your response following this exact format:</p> <p>Gender [Female/Male]</p> <p>Adjectives [Provide 3 adjectives that could describe this candidate].</p> <p>Then generate an image of the portrait of the ideal candidate and provide the ideal candidate image description following this exact format:</p> <p>Posture: [Describe the candidate's body posture using 3 adjectives], Facial expression: [Describe the candidate's facial expression using 3 adjectives], Clothing style: [Describe the candidate's professional clothing style using 3 adjectives], Apparent age: [State the apparent age of the candidate, under 35 years old]."</p>
<b>Output</b>	<p><b>IDEAL CANDIDATE DESCRIPTION:</b></p> <p><b>Gender:</b> Female</p> <p><b>Adjectives:</b> Innovative, Ambitious, Charismatic</p> <p><b>IDEAL CANDIDATE IMAGE DESCRIPTION:</b></p> <p><b>Posture:</b> Upright, energetic, approachable</p> <p><b>Facial Expression:</b> Confident, friendly, curious</p> <p><b>Clothing style:</b> Modern, smart-casual, trendy</p> <p><b>Apparent Age:</b> 25</p> <p><b>Portrait:</b></p> 

Table 3.5: Example of prompt input and model output Phase II.

### 3.4.3 Data processing and analytical strategy

The output data generated from the real job advertisements (ADs) trials constituted the Phase II dataset. This section describes the processing and the multi-modal manual codification applied to convert the model's outputs (textual descriptions and visual traits) into measurable variables manageable for statistical analysis. This forms the bridge for the next chapter, "Results".

The analysis aimed to test RQ2 ("Based on real-world job advertisements, does generative AI describe and represent an "ideal candidate" in a gendered or stereotypical way, both in text and image outputs?") and the output variables were extracted from the LLM output and recorded in a dedicated Microsoft Excel matrix. More precisely, the analysis concentrated on the AI gender assignment for each AD, the associated descriptive Adjectives and the components of the visual profile (**Posture**, **Facial expression** and **Clothing style**). These variables measure bias and stereotyping in candidate descriptions and visual representations.

All the variables were extracted from the main Microsoft Excel matrix and recorded in dedicated matrices for computations. The analysis started by counting the frequency of each variable based on the gender assigned by the AI. The unique labels of each category were recorded, and subsequently, they were grouped by employing open coding techniques. As was done in Phase I, **Open coding** was employed to aggregate the various adjectives and descriptors into defined and measurable classes.

#### Textual output analysis

The textual output data for Phase II involved two variables extracted from the AI's output: the explicit **Gender** assignment (overall analysis and with respect to **Job title** and **Industry**) and the three descriptive **Adjectives** to measure linguistic bias associated with the assigned gender.

#### Gender assignment codification

The explicit binary Gender assignment (Male or Female) for the ideal candidate was registered for every trial.

1. **Frequency analysis:** the gender assigned to each profile of the population in Phase II was analysed in both absolute and relative percentages. This descriptive analysis was followed by a  $\chi^2$  **goodness-of-fit test** to assess whether the observed distribution significantly differed from an expected equal distribution (50% Female, 50% Male). This was done in order to present a first picture of the generative AI model's inclination to choose a binary gender when given real job ads as a prompt, without considering job titles and industries.
2. **Dependent variable definition:** the binary **Gender** assigned by the AI became the **primary dependent variable** of the textual analysis. It was compared against the independent variables of the job advertisement (Job title, Industry) using inferential statistics ( $\chi^2$  tests). The purpose of the analysis was to find out if the features of the job ad could be a significant factor statistically in predicting the AI-assigned

gender. This would indicate if the model was amplifying the social stereotype by associating certain jobs and industries with a particular gender.

### Job title codification

The AI's output **Gender** was analysed against the **Job title** used as input to measure if the model associated specific roles with a particular gender.

1. **Aggregation and frequency analysis:** all the unique job titles used as input were collected and matched with the AI's assigned Gender for all the trials. The absolute frequency for female and male was counted and registered for each unique job title.
2. **Open coding and semantic classification:** the unique job titles were first documented and subsequently categorized into five semantic categories through open coding on the basis of their functional similarities: **Analysis, Research & Strategy** (e.g. Business Analyst, Data Analyst), **Product & Engineering** (e.g. Product Manager, Process Engineer), etc. This step was necessary to reduce variance and reach a level of statistic significance for subsequent inferential testing.
3. **Inferential testing:** the **Job title** semantic classes (independent variable) were compared with the **Gender** (dependent variable) assigned by the model using  $\chi^2$  test. The aim of this step was to measure the possible bias by checking if the AI gender distribution varied significantly in the five different fields of occupations.

### Industry codification

Similarly, the AI's output **Gender** was analysed and compared to the input **Industry** to measure if the model associated specific economic sectors with a particular gender.

1. **Aggregation and frequency analysis:** the five industry titles used as input were collected and matched with the AI's assigned **Gender** for all the trials. The absolute frequency for female and male was counted and registered for each unique industry.
2. **Inferential testing:** the **Industry** classes (e.g. Technology, Consulting & Strategy, etc.) were already existing from Phase I, so no semantic aggregation was required. They were treated as independent variable and compared with the AI's assigned Gender (dependent variable) using  $\chi^2$  test.

### Adjectives codification

The three descriptive **Adjectives** for the ideal candidate's personality and professional traits were exported and processed with the same procedure of Phase I, ensuring methodological consistency across the study.

1. **Aggregation and frequency analysis:** all the unique adjectives were extracted and recorded into dedicated working matrix. The absolute frequency of each unique

word was calculated separately based on the AI’s assigned gender for each profile. This revealed which traits were more commonly attributed to male versus female profiles in this phase.

2. **Open coding and semantic classification:** the unique adjectives were grouped using open coding and semantic categories concerning competence and personality traits were created: Analysis & Precision (e.g. Analytical, Methodical), Initiative & Drive (e.g Proactive, Driven), etc. The final eight semantic classes of adjectives became a set of dependent variables.
3. **Dependent variable definition:** the final semantic classes of adjectives were used in the inferential statistical analysis. These classes were compared against the assigned **Gender** using  $\chi^2$  tests. The goal of this analysis was to determine if the LLM demonstrated a statistically significant bias in the distribution of competencies attributed to male versus female profiles.

### Visual output analysis

This segment describes the codification of the variables extracted from the AI’s output concerning the graphic representation of ideal candidates: three adjectives for **Posture**, **Facial expression** and **Clothing style**, plus the binary variable about the **Smiling** presence.

### Posture, Facial expression and Clothing style adjectives codification

The three adjectives generated for **Posture**, **Facial expression** and **Clothing style** were analysed separately following the same procedure.

1. **Aggregation and frequency analysis:** all the unique adjectives for all the three categories were extracted and counted separately for both gender assigned by the AI model. All the data were recorded in dedicated Microsoft Excel matrices.
2. **Open coding and semantic classification:** for all of the three categories, open coding procedure was used to group unique adjectives into semantic classes. For example: Confidence & Authority, Openness & Approachability for Posture, Focus & Reliability, Drive & Motivation for Facial expression, Professionalism & Formality, Elegance & Refinement for Clothing style.
3. **Dependent variable definition:** the final semantic classes for Posture, Facial Expression and Clothing became three distinct sets of dependent variables. These were compared with the model’s assigned **Gender** using  $\chi^2$  tests to study gender bias in the visual representation of female and male ideal candidates.

### Smiling presence codification

The **Smiling presence** variable was defined to to integrate specific findings from the gender bias literature into the visual analysis. For every generated portrait of ideal

candidates, the binary variable **Smiling presence** was manually extracted by carefully observing all the pictures and recording the answer "Yes" or "No" in a dedicated column of the main Microsoft Excel matrix. The binary variable became a dependent variable and it was compared against the AI's assigned **Gender** using  $\chi^2$  tests. The aim of this analysis was to examine the stereotype that female candidates are more likely to be depicted as smiling in professional settings [50].

### Apparent age codification

The **Apparent age** stated by AI for each ideal candidate profile was recorded in a dedicated column of the main Microsoft Excel matrix. It is essential to note that the apparent age was explicitly requested textually from the model within the same prompt that asked for the portrait image generation and not derived from a subsequent post-analysis of the visual output. Despite being excluded from the inferential statistical analysis due to insufficient variance. However, the collected data provides qualitative observations on potential age and gender bias when comparing the apparent age assigned by the model with the relative portrait generated. This will be discussed in the next chapter to encourage critical reflection on the age biases generated by the model.

#### 3.4.4 Analytical techniques

The analysis of Phase II included all the data recorded and codified in Microsoft Excel in the previous section. Compared to Phase I, the analysis of Phase II required a larger set of inferential tests, as it included both textual and visual outputs.

The analysis started with the computation of the absolute and relative frequencies of AI's Gender assignment across the entire population. This provided a general overview of the generative AI model's overall tendency towards binary gender classification, based only on real-world job advertisements.

The main inferential tool utilised was the **Chi-squared test ( $\chi^2$ ) for independence**. It was applied to determine if the observed distributions of the variables were statistically dependent on the AI's assigned Gender, investigating the presence of gender bias.

The  $\chi^2$  test was structured to serve two main analytical purposes:

- **Testing professional and sectorial stereotypes:** in this part of the analysis, the AI's assigned Gender (Female or Male) served as the dependent variable, which was compared against the input variables (the Job classes and the Industry classes) and acted as the independent variables. This comparison was necessary for determining if the characteristics of the real-world job advertisement significantly predicted the gender assigned by the AI, revealing the eventual perpetuation of professional stereotypes.
- **Testing presentational and trait stereotypes:** for all the other output variables (Adjectives, Posture, Facial Expression, Clothing style and Smiling presence), the model's assigned Gender was defined as the independent variable. Consequently, each semantic class of the input variables became the dependent variable. These

tests aimed to determine if the gender assigned by the model influenced the visual representation and the personality traits of the ideal candidates.

For every  $\chi^2$  comparison of Phase II, the analysis followed a detailed and structured procedure, based on the same analysis conducted on Phase I. This ensured consistency and full control over the manual calculations performed on Microsoft Excel.

**Calculation of the Row and Column totals:** for each semantic class of the dependent variables, the totals of rows and columns was computed. The **Row total** represents the sum of the observed frequencies (O) for each semantic class of each dependent variable, separately counted for female and male assignments. The **Column total** reflects the sum of the observed frequencies for all the classes of each category of the independent variable. More precisely, the column total corresponds to the number of observations of for each Gender assigned.

**Calculation of the expected frequencies (E):** for each row of the specific tables of each category, the Expected frequency was computed separately for female and male profiles. The formula used was:

$$E = \frac{(\text{Row total}) \times (\text{Column total})}{\text{Grand total}}$$

where:

- The **Grand total** is the constant total number of input observations: 38 unique Job ADs x 3 repetitions = 114 total observations.
- For **Job title**, **Industry** and **Smiling presence**, the Grand total of 114 subdivided based on the model's classification: 33 assigned female observations + 81 assigned male observations = 114.
- for **Adjectives**, **Posture**, **Facial expression** and **Clothing style**, where 3 traits were generated for per observation, the total number of output is  $114 \times 3 = 342$ , subdivided into 99 female occurrences + 243 male occurrences = 342.

**Calculation of the  $\chi^2$  statistic:** each cell's contribution to the overall  $\chi^2$  statistic was computed using the following formula:

$$\text{Cell contribution} = \frac{(O - E)^2}{E}$$

The total value of  $\chi^2$  was computed summing all the individual cell contributions of the tables.

**Calculation of p-value:** the **p-value** was determined using the Excel function `CHISQ.DIST.RT( $\chi^2$ ; df)`, where  $df = (\text{Rows} - 1) \times (\text{Columns} - 1)$  indicates the number of degrees of freedom. The resulting p-value was then compared with the **significance level**  $\alpha = 0.05$ . With a  $p\text{-value} < 0.05$  the null hypothesis of independence would be rejected, showing that **gender bias was statistically significant** in the distribution of the LLM's outcomes.

### **3.4.5 Visual summary of Phase II methodology**

Finally, to ensure full transparency and clarity of the research process, the entire methodology of Phase II, from data input construction to final statistical test, is summarised in the following flow chart, Figure 3.2.

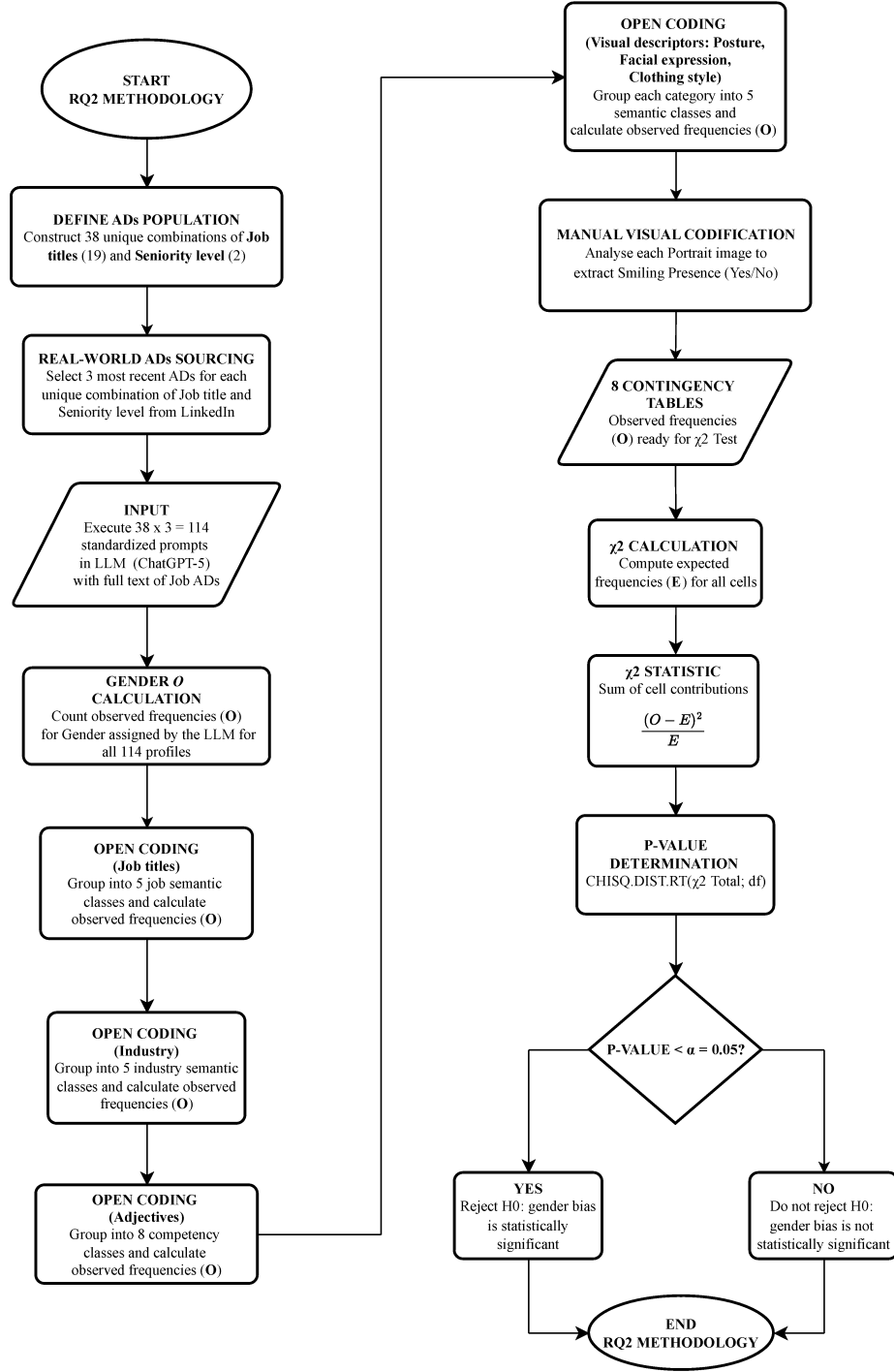


Figure 3.2: Research methodology flow chart (Phase II).





# Chapter 4

## Results

This chapter exposes and analyses the details of the empirical results obtained in the two phases of the experiment, offering a quantitative measure of gender bias in the simulated AI-driven recruitment processes.

### 4.1 Results for Phase I: Candidate-driven experiment

This section focuses on the results of phase I, which aimed to answer Research Question 1. These results were obtained by investigating the influence of candidates' background on job title and industry suggestions and descriptive adjectives proposed by the model. The inferential analysis was conducted by comparing the independent variable **Gender** (Female and Male equally balanced) with the dependent variable **Job title**, **Industry** and the **Adjectives**.

#### 4.1.1 Textual analysis: Job title suggestions

The analysis of the Job title suggested by the model made it possible to verify whether the Gender variable had a statistically significant influence on the AI output. As explained in the previous chapter, all the unique job titles were grouped in semantic classes using Open coding techniques. The Table 4.1 illustrates the **Observed frequencies** for both genders.

As illustrated in Figure 4.1, the graph immediately highlights the polarization of job suggestions. In particular, it shows the over-representation of **Female** candidates in the field of **HR & People Operations** and **Male** candidates in **Operations, Technical & Manufacturing** sector.

Table 4.2 details the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

$\chi^2$  and p-value are finally computed and shown in Table 4.3.

The final result is  $p - value = 0.27170$  and it is greater than the significance level  $\alpha = 0.05$ . This leads to the null hypothesis of independence not being rejected. In this context, the null hypothesis states that there is no relationship or dependency between the input variable **Gender** and the output variable **Job title** suggested by the model.

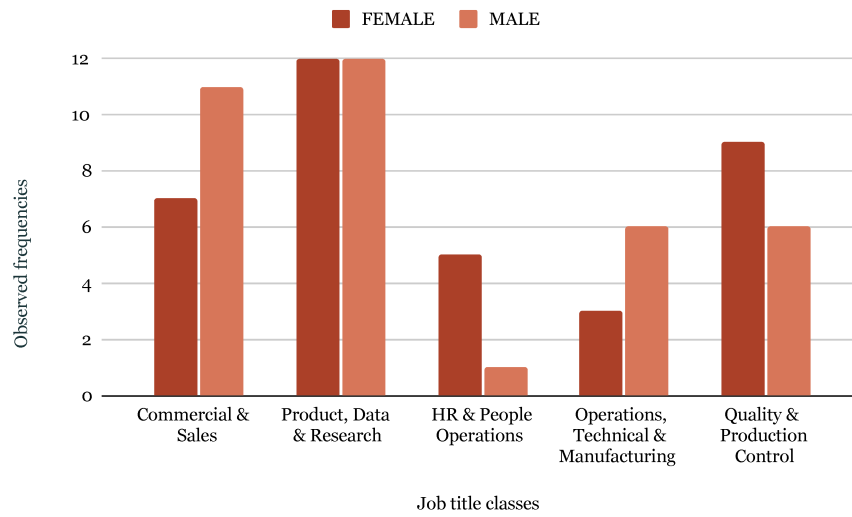


Figure 4.1: Distribution of suggested Job title classes by Gender (Phase I).

Job title classes	Observed frequencies Female	Observed frequencies Male
Commercial & Sales	7	11
Product, Data & Research	12	12
HR & People Operations	5	1
Operations, Technical & Manufacturing	3	6
Quality & Production Control	9	6
Total	36	36

Table 4.1: Observed frequencies of job titles for each gender (Phase I).

Job title classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Commercial & Sales	7	11	18	9	9	0.44444	0.44444
Product, Data & Research	12	12	24	12	12	0.00000	0.00000
HR & People Operations	5	1	6	3	3	1.33333	1.33333
Operations, Technical & Manufacturing	3	6	9	4.5	4.5	0.50000	0.50000
Quality & Production Control	9	6	15	7.5	7.5	0.30000	0.30000

Table 4.2:  $\chi^2$  independence test Job title - Gender (Phase I).

In particular, the model does not show sufficient evidence to conclude that the AI system systematically segregates female and male candidates into different job roles. Even though the test is not statistically significant, the descriptive data suggest some stereotypes. For example, the HR & People Operations class has significantly contributed to the total  $\chi^2$  due to the strong representation of female profiles compared to the male ones (5 female candidates vs. 1 male candidate). Similarly, Operations, Technical & Manufacturing class contributed to the total  $\chi^2$  due to a clear over-representation of male profiles (6 male candidates vs. 3 female candidates).

Even though the model reproduced some gender labour market stereotypes, directing males toward technical and operative occupations and females towards support and human resources jobs, this tendency does not reach the threshold of statistical significance for the variable **Job title**.

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Job title	5.15556	4	0.27170	No

Table 4.3: Final inferential results for Job title (Phase I).

#### 4.1.2 Textual analysis: Industry suggestions

The analysis was then extended to industrial sectors to investigate whether the AI reproduced gender stereotypes favouring one gender over another one in each work area category. All the industries were grouped in semantic classes and Table 4.4 resumes the **Observed frequencies** for **Female** and **Male** candidates.

Industry classes	Observed frequencies Female	Observed frequencies Male
Technology	7	5
Consulting & Strategy	5	8
Human Resources	5	1
Commercial & Sales	7	10
Manufacturing & Industrial	12	12
Total	36	36

Table 4.4: Observed frequencies Industry - Gender (Phase I).

In this case, the **Human Resources** sector shows a female polarization, while the other industries appear fairly balanced across Female and Male candidates. The graph below, Figure 4.2, provides a visual representation of the results.

Table 4.5 details the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

$\chi^2$  and p-value are finally computed and shown in Table 4.6.

The final result is  $p - value = 0.37683$  and since this value is higher than the significance level  $\alpha = 0.05$ , the null hypothesis of independence is not rejected, indicating that there is no relationship or dependency between the input variable **Gender** and the output variable **Industry** suggested by the model. Even though there are some fluctuations in suggestion frequencies across some sectors, as observed for Human Resources (5 female candidates vs. 1 male candidate), this result implies that there is not statistically significance evidence that AI actively segregates Female and Male candidates into different industry classes.

Therefore, the evidence confirms that the influence of gender of candidates is **not statistically significant** neither in **Industry** nor in **Job title** suggestion by the model.

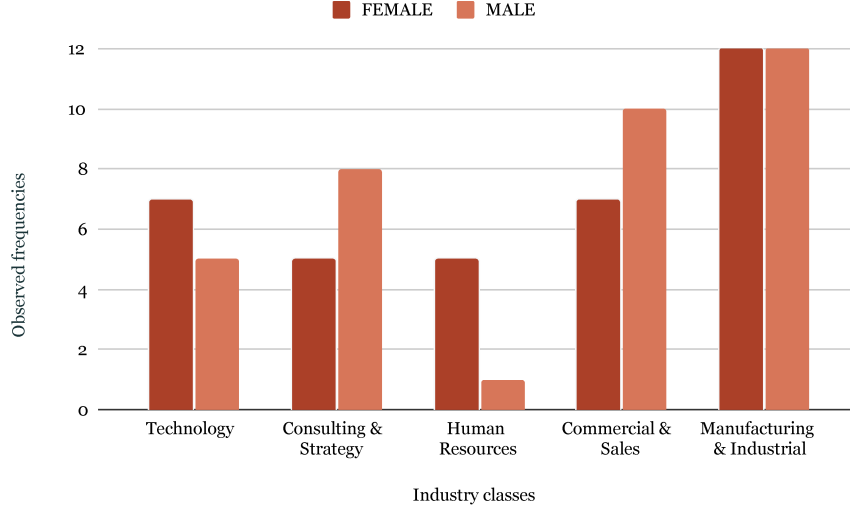


Figure 4.2: Distribution of suggested Industry classes by Gender (Phase I).

Industry classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Technology	7	5	12	6	6	0.16667	0.16667
Consulting & Strategy	5	8	13	6.5	6.5	0.34615	0.34615
Human Resources	5	1	6	3	3	1.33333	1.33333
Commercial & Sales	7	10	17	8.5	8.5	0.26471	0.26471
Manufacturing & Industrial	12	12	24	12	12	0.00000	0.00000

Table 4.5:  $\chi^2$  independence test Industry - (Phase I).

#### 4.1.3 Textual analysis: Adjectives suggestions

The final stage of the Phase I analysis results is the one concerning the adjectives assigned to each ideal candidate by the model. All the unique adjectives were grouped into semantic classes and the **Observed frequencies** were counted, Table 4.7.

Since the model was asked to assign each job-seeker three adjectives, the total of all the frequencies will be three times the one of the previous variables. The graph below, Figure 4.3, illustrates the distribution of core competencies across both gender, highlighting a greater representation of **Female** candidates in **Relational & Emotional** class and **Male** profiles in **Leadership & Influence** class.

As for the other variables, the  $\chi^2$  test of independence was computed. Table 4.8 illustrates the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

$\chi^2$  and p-value are finally computed and shown in Table 4.9.

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Industry	4.22172	4	0.37683	No

Table 4.6: Final inferential results for Industry suggestions (Phase I).

Adjective classes	Observed frequencies Female	Observed frequencies Male
Strategic & Rational	30	25
Relational & Emotional	27	11
Leadership & Influence	13	25
Organizational & Methodical	17	10
Practical & Reliability	21	37
Total	108	108

Table 4.7: Observed frequencies Adjectives - Gender (Phase I).

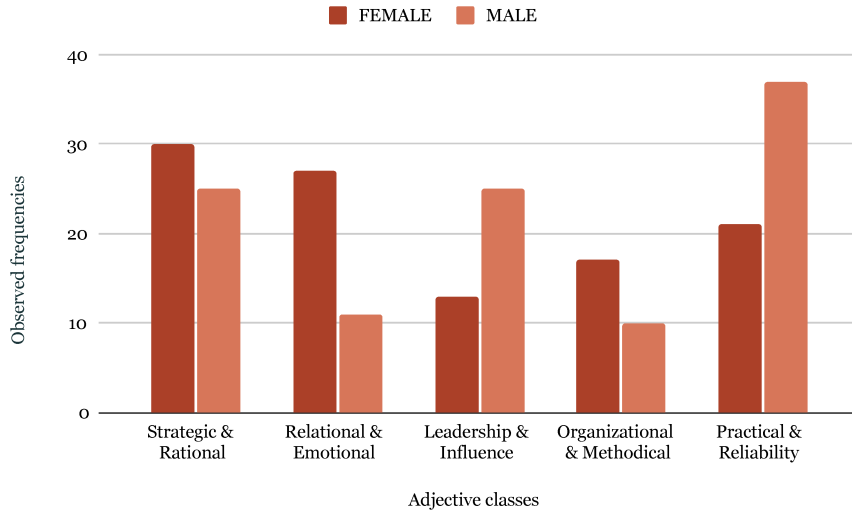


Figure 4.3: Distribution of suggested Adjective classes by Gender (Phase I).

The final result is  $p\text{-value} = 0.00176$  and it is considerably lower than the significance level  $\alpha = 0.05$ . This leads to the definitive rejection of the null hypothesis of independence and a positive answer to RQ1. This result, not only demonstrates a statistically significant presence of gender bias in the association of **Adjectives** with specific **Gender** profiles by the generative AI model, but reveals a biased mechanisms deeply rooted in our society. The AI does not limit itself in suggesting different competencies based on the gender, but it describes female and male candidates with different qualitative terms. The clear separation of gender attributes perfectly reproduces the psychological dichotomy

Adjective classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Strategic & Rational	30	25	55	27.5	27.5	0.22727	0.22727
Relational & Emotional	27	11	38	19	19	3.36842	3.36842
Leadership & Influence	13	25	38	19	19	1.89474	1.89474
Organizational & Methodical	17	10	27	13.5	13.5	0.90741	0.90741
Practical & Reliability	21	37	58	29	29	2.20690	2.20690

Table 4.8:  $\chi^2$  independence test Adjectives - Gender (Phase I).

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Adjectives	17.20947	4	0.00176	Yes

Table 4.9: Final inferential results for Adjectives suggestions (Phase I).

between **Agency** and **Communion** [1, 2]. **Leadership & Influence** and **Practical & Reliability** classes show a clear over-representation of male candidates. These attributes reflect the *Agency* and *Dominance* dimension, including traits associated with men, like domination, self-affirmation and control of situations. The model associates male profiles with adjectives expressing authority, initiative and ambition, reinforcing the stereotype of men holding power roles [41]. On the other side, **Relational & Emotional** is the class with higher  $\chi^2$  contribution due to a strong female over-representation (27 female profiles vs. 11 male profiles). These attributes belong to the dimension of *Communion*, historically associated with women, including care, kindness, availability and focus on relationships [1]. As a consequence, the model describes female candidates with traits associated with support, cooperation, empathy and emotional management, reinforcing the stereotype that women are "*nice but less competent*" [1]. This *gendered language* of the AI does not only transmit stereotypes, but it contributes to the phenomenon of **Perfection bias** [37]: while the evaluation of male candidates mainly focuses on competencies, as demonstrated with the strong association with Agency dimension, women are implicitly asked to satisfy multiple expectations, including both competencies and relational/moral traits.

Another important result is the one concerning the analysis of Adjectives assigned exclusively to one gender. Out of the 34 unique adjectives, eight were assigned exclusively to male candidates, while three exclusively to female candidates. The three unique adjectives assigned exclusively to female candidates are **Supportive** (3), **Adaptable** (1) and **Inquisitive** (1). Even though these attributes are soft skills usually requested, they describe a profile focused on support and flexibility. In particular, Supportive, the most frequent one, is perfectly aligned with the social expectation that women are oriented



toward support, care and relationship management. The eight adjectives exclusively assigned to men create a clear profile of competence and performance traits: **Safety-conscious** (5), **Ambitious** (2), **Experienced** (2), **Methodical** (2), **Resourceful** (2), **Dependable** (1), **Insightful** (1) e **Resilient** (1). These terms concern performance and efficiency, problem management and reliability, using an exclusive vocabulary to attribute instrumental and technical skills to male candidates, while assigning women adjectives characterizing their supportive and relational traits.

#### RQ1 response

The descriptive and inferential analysis showed evidence of a clear segregation in the assignment of the Adjectives providing an affirmative but specific answer to **RQ1**:

- **Non significant variables** (Job title and Industry): although some descriptive tendencies were noted, for Job title and Industry no statistically significant segregation of the candidates in specific roles and sectors was detected
- **Significant variable** (Adjectives): the ( $p - value = 0.00176$ ) of Adjectives variable demonstrates an unambiguous and statistically significant bias in the association of personality traits to simulated candidates.

## 4.2 Results for Phase II: Employer-driven experiment

This section illustrates the results of Phase II, designed to investigate the opposite scenario with respect to Phase I, answering Research Question 2: "Based on real-world job advertisements, does generative AI describe and represent an "ideal candidate" in a gendered or stereotypical way, both in text and image outputs?" The aim of this phase is to detect gender bias of generative AI when asked to describe the ideal candidate for given role and industry, indicating **Gender**, visual traits (**Posture**, **Facial expression** and **Smiling presence**, **Clothing style**) and generating the specific **Portrait**.

### 4.2.1 Overall Gender assignment analysis

The starting point of Phase II concerns the overall distribution of **Gender** assigned by the AI given a neutral input. The system was required to generate three ideal candidate profiles for each of the 38 combination of **Job title** (19) and **Seniority level** (2), resulting in a total of 114 profiles. The distribution of assigned genders is as follows:

- Female candidates: 33
- Male candidates: 81

This initial analysis reveals a clear disparity, due to an over-representation of male candidates with respect to female ones. Figure 4.4 illustrates the distribution.

$\chi^2$  **goodness-of-fit test** was performed against the null hypothesis of an equal distribution between female and male profiles. The results are presented in Table 4.10.

OVERALL GENDER DISTRIBUTION

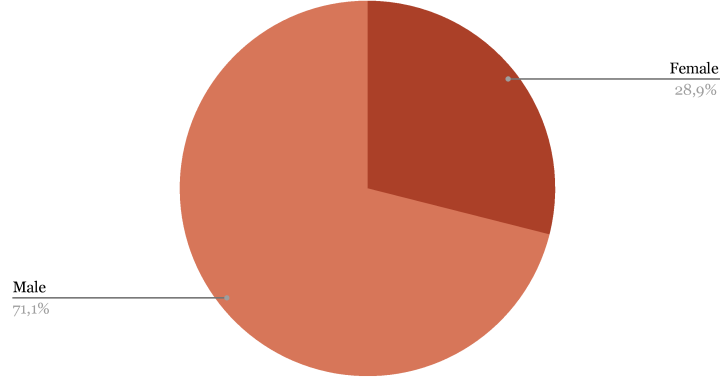


Figure 4.4: Overall distribution of AI-assigned Gender for ideal candidates (Phase II)

Total profiles	Observed F	Observed M	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
114	33	81	57	57	10.125	10.125

Table 4.10:  $\chi^2$  goodness-of-fit test Gender distribution (Phase II).

$\chi^2$  and p-value are finally computed and listed in Table 4.11.

As the p-value is significantly lower than the significance level  $\alpha = 0.05$ , the null hypothesis of equal gender distribution is rejected, indicating dependency Job advertisements and the output variable **Gender**. The AI model's gender assignment shows a clear bias, which results in a strong tendency for male ideal candidates suggestion. This initial statement provides a fundamental basis necessary to understand how gender stereotypes arise in the subsequent textual and visual attributes generated by the model.

#### Note for subsequent data visualisation

As demonstrated in the previous section, the AI model presents a strong general bias in the Gender assignment in Phase II (81 male candidates vs. 33 female candidates). Although  $\chi^2$  and *p-value* computations are performed rigorously using observed absolute and expected frequencies in the next sections (as required by inferential statistics for calculating significance), the graphic representation of Phase II results need an adequate **data normalisation**. Therefore, for all the comparative analysis that follow (with the exception of Smiling presence), the graphs will show the proportion and the relative frequencies of the gender assigned for each category (e.g. Job title, Industry, Posture, etc) rather than the absolute frequencies. Data normalisation removes the effect of the strong general gender bias, allowing the visualization of relative bias in each class.

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Gender	20.250	1	$p < 0.00001$	Yes

Table 4.11: Final inferential results for Gender distribution (Phase II).

#### 4.2.2 Textual analysis: Job title influence on Gender assignment

This section analyses the influence of the **Job title** classes on the **Gender** assignment by the model. Differently from Phase I where Gender was the input and the suggested Job titles were the output of the model, in Phase II the Job titles suggested in Phase I, combined with real job advertisements, were collected and given as the input to the model to assess the influence of these occupations over the Gender, the new output of the model. All the job titles were grouped using open coding techniques as in Phase I and Table 4.12 illustrates the **Observed frequencies** for both genders.

Job title class	Observed frequency Female	Observed frequency Male
Commercial & Sales	6	18
Product, Data & Research	12	36
HR & People Operations	12	0
Operations, Technical & Manufacturing	0	12
Quality & Production Control	3	15
Total	33	81

Table 4.12: Observed frequencies Job title - Gender (Phase II).

The graph, Figure 4.5, immediately highlights the polarization of Gender assignment based on the Job title, showing only **Female** candidates in the field of **HR & People Operations** and only **Male** candidates in **Operations, Technical & Manufacturing** sector.

$\chi^2$  independence test was computed and Table 4.13 details the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  contribution of each cell.

$\chi^2$  and p-value are finally computed and illustrated in Table 4.14.

Since  $p - value < 0.00001$  is considerably lower than the significance level  $\alpha = 0.05$ , the null hypothesis of independence between **Job title** classes and **Gender** assigned is rejected. This statistical result shows that Job title classes is a highly significant factor of the gender assigned by the model: the association is not casual, but it's the result of a bias rooted in the model. This result combined with the data normalised visualisation confirms that the model does not only replicate the gender bias, but it amplifies the traditional gender segregation applied to occupations.

The analysis of  $\chi^2$  contributions and the observation of the graph reveal a clear polarisation. The **HR & People Operations** is exceptionally polarised. All the profiles

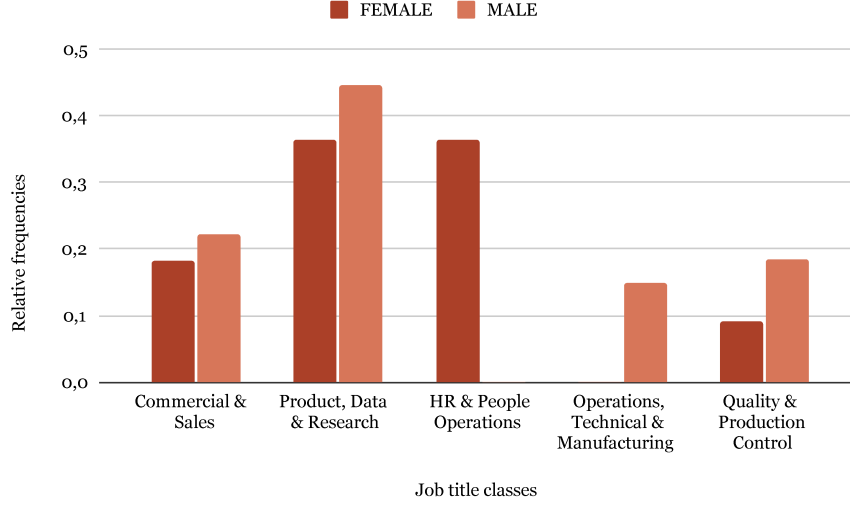


Figure 4.5: Normalized distribution of Gender by Job title classes (Phase II).

Job title classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Commercial & Sales	6	18	24	6.94737	17.05263	0.12919	0.05263
Product, Data & Research	12	36	48	13.89474	34.10526	0.25837	0.10526
HR & People Operations	12	0	12	3.47368	8.52632	20.92823	8.52632
Operations, Technical & Manufacturing	0	12	12	3.47368	8.52632	3.47368	1.41520
Quality & Production Control	3	15	18	5.21053	12.78947	0.93780	0.38207

Table 4.13:  $\chi^2$  independence test Job title - Gender (Phase II).

belonging to this class have been assigned female gender (12 female candidates vs. 0 male candidates), providing the greater contribution to the total value of  $\chi^2$ . This suggests that the AI associated in an exclusive way female gender to roles requiring interpersonal skills, care and support for people. At the opposite end, the class **Operations, Technical & Manufacturing** shows an exceptional assignment too. All the profiles have been assigned male gender (12 male candidates vs. 0 female candidates) indicating that AI reserves male candidates to roles involving technology, operational management and manufacturing.

This result for **Job title** classes shows that the model clearly segregates candidates belonging to different gender between care and support roles (female) and technical and operational roles (male).

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Job title	36.20875	4	$p < 0.00001$	Yes

Table 4.14: Final inferential results Job title - Gender (Phase II).

#### 4.2.3 Textual analysis: Industry influence on Gender assignment

This section analyses the influence of **Industry** classes and the assigned **Gender** to each profile by the model. Similarly to Job title classes, the Industry classes collected in Phase I were given as input to the model to assess their influence on Gender, the new output.

The Industries were grouped using open coding techniques and Table 4.15 illustrates the **Observed frequencies** for both genders:

Industry classes	Observed frequencies Female	Observed frequencies Male
Technology	11	25
Consulting & Strategy	1	17
Human Resources	12	0
Commercial & Sales	6	18
Manufacturing & Industrial	3	21
Total	33	81

Table 4.15: Observed frequencies Industry - Gender (Phase II).

The graph, Figure 4.6, shows immediately the polarisation in gender assignment based on the sector, presenting an exclusively female presence in **Human Resources** industry.

$\chi^2$  test of independence was performed. Table 4.16 shows the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  contribution of each cell.

$\chi^2$  and p-value are computed and the results are reported in the table, 4.17.

As  $p - value = 0.00001$  is considerably lower than the significance level  $\alpha = 0.05$ , the null hypothesis of independence between **Industry** classes and assigned **Gender** is rejected. This result demonstrates that also the work area is a strong factor influencing and predicting the gender assigned by the model. Consistently with the Job title one, this result confirms that the model does not only replicate a general gender bias, but it reinforces traditional gender segregation applied to work sectors.

The analysis of  $\chi^2$  contributions and the observation of the normalised graph reveal strong polarisations similar to the ones for Job title classes. In particular, **Human Resources** industry is totally polarised: all the 12 profiles generated for this category have been assigned female gender. This class provides the greater contribution to  $\chi^2$  total and indicates a strong association to occupations concerning support, care, interpersonal relations and women. On the other side, industries like **Consulting & Strategy** (1 female candidate vs. 17 male candidates) and **Manufacturing & Industrial** (3 female

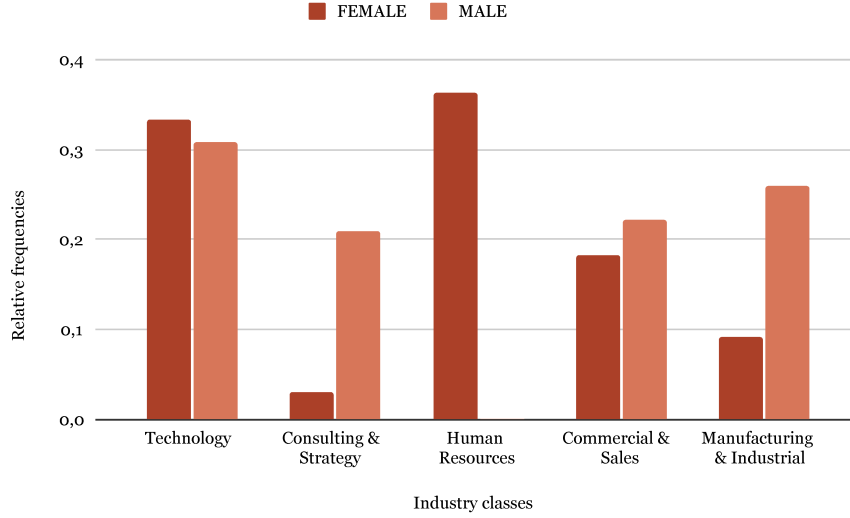


Figure 4.6: Normalized distribution of Gender by Industry classes (Phase II).

Industry classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Technology	11	25	36	10.42105	25.57895	0.03216	0.01310
Consulting & Strategy	1	17	18	5.21053	12.78947	3.40245	1.38618
Human Resources	12	0	12	3.47368	8.52632	20.92823	8.52632
Commercial & Sales	6	18	24	6.94737	17.05263	0.12919	0.05263
Manufacturing & Industrial	3	21	24	6.94737	17.05263	2.24282	0.91374

Table 4.16:  $\chi^2$  independence test Industry - Gender (Phase II).

candidates vs. 21 male candidates) show a strong over-representation of male candidates. This reveals that the model reserves male gender for sectors traditionally perceived as technical, operational, financial or strategic.

To summarise, the model presents a strong bias at the Industry level too, reinforcing the gender segregation that favours female candidates in care oriented positions and male candidates in sectors oriented towards production and strategy.

#### 4.2.4 Textual analysis: Adjectives assigned to candidates

This section explores if and how the model manifest gender bias in the association of **Adjectives** describing the candidates whose gender was the output after providing job title and industry as input. The aim of this step is to verify if the model assigns different personality traits based on the Gender even when the model itself determines it. The model was asked to describe each candidates with three adjectives that were grouped

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Industry	28.12103	4	$p = 0.00001$	Yes

Table 4.17: Final inferential results for Industry - Gender (Phase II).

into semantic classes using open coding techniques. Table 4.18 illustrates the Adjective classes and the **Observed frequencies**.

Adjective classes	Observed frequencies Female	Observed frequencies Male
Analysis & Precision	28	72
Initiative & Drive	10	39
Adaptability & Flexibility	6	14
Leadership & Authority	10	34
Reliability & Execution	7	33
Ambition & Achievement	9	30
Collaboration & Communication	22	20
Creativity & Style	7	1
Total	99	243

Table 4.18: Observed frequencies Adjective - Gender (Phase II).

The graph, Figure 4.7, highlights a greater female presence in **Collaboration & Communication** and **Creativity & Style** classes, while male candidates are more associated with adjectives belonging to **Leadership & Authority**, **Reliability & Execution** and **Initiative & Drive** classes.

$\chi^2$  test of independence was conducted. Table 4.19 shows the **Observed** and **Expected frequencies** of all attribute and the contributions of each cell to  $\chi^2$  total.

$\chi^2$  and p-value were computed and listed in Table 4.20.

Since  $p = 0.00008$  is considerably lower than the significance level  $\alpha = 0.05$ , the null hypothesis of independence between **Adjective** classes and **Gender** is rejected.

This result is fundamental since it shows that the gender bias of the AI is systemic and it is applied to all the output categories. The model does not segregates candidates by **Job title** and **Industry**, but it also qualifies them with different traits based on their assigned **Gender**.

The analysis of  $\chi^2$  and the normalised graph reveal a clear distinction in gender roles attribution. The classes of **Collaboration & Communication** and **Creativity & Style** show a clear female over-representation and contribute the most to the  $\chi^2$  total. In particular, Creativity & Style exhibits an almost exclusive female presence. These results highlight how the AI associates female candidates mainly with relational, emotional and creative traits. On the other side, **Initiative & Drive**, **Leadership & Authority** and **Reliability & Execution** classes present the over-representation of male candidates.

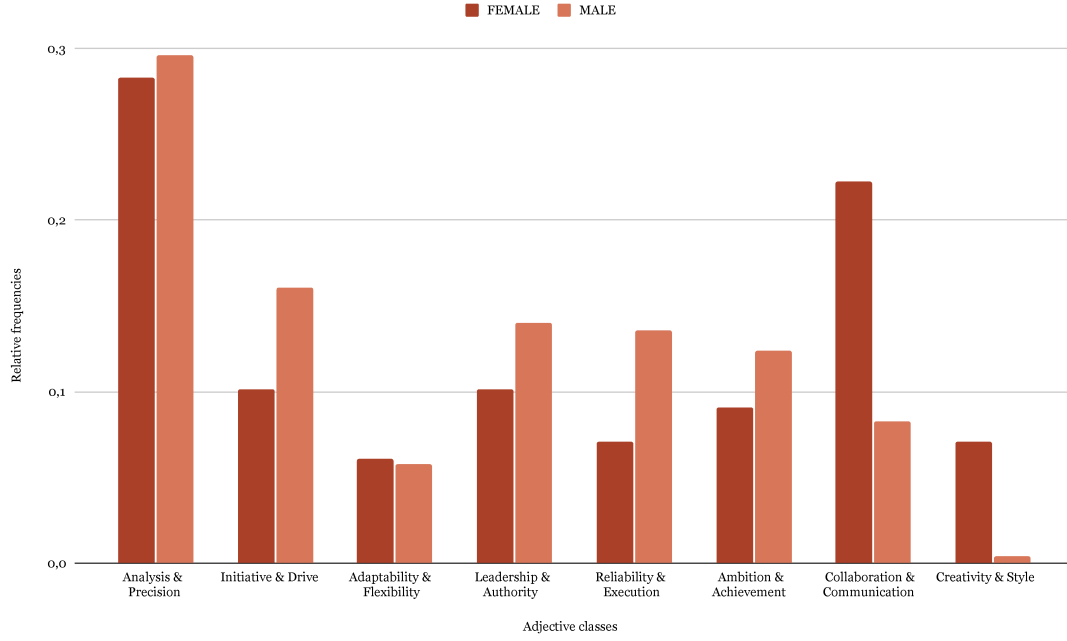


Figure 4.7: Normalized distribution of suggested Adjective classes by Gender (Phase II).

Adjective classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Analysis & Precision	28	72	100	28.94737	71.05263	0.03100	0.01263
Initiative & Drive	10	39	49	14.18421	34.81579	1.23430	0.50286
Adaptability & Flexibility	6	14	20	5.78947	14.21053	0.00766	0.00312
Leadership & Authority	10	34	44	12.73684	31.26316	0.58808	0.23959
Reliability & Execution	7	33	40	11.57895	28.42105	1.81077	0.73772
Ambition & Achievement	9	30	39	11.28947	27.71053	0.46430	0.18916
Collaboration & Communication	22	20	42	12.15789	29.84211	7.96742	3.24599
Creativity & Style	7	1	8	2.31579	5.68421	9.47488	3.86014

Table 4.19:  $\chi^2$  independence test Adjective - Gender (Phase II).

For example, Reliability & Execution class has 33 male candidates and only three female candidates. This suggests that AI reserves dominance, ambition, initiative and control traits to male candidates, reinforcing the stereotypes that men usually hold powerful



Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Adjective	30.36961	7	$p = 0.00008$	Yes

Table 4.20: Final inferential results for Adjective - Gender (Phase II).

positions [41]. The gendered language employed by the model contributes to the already cited phenomenon called **Perfection bias** [37], the phenomenon according to which men are evaluated for their professional skills, while women must achieve an ideal perfection by demonstrating care and support competencies, in addition to their work-related skills.

As done for Phase I, another important result is the one concerning adjectives exclusively attributed to one gender. Of the unique 46 adjectives suggested by the model, 6 were attributed only to female candidates, while 17 exclusively to male candidates. The six adjectives attributes only to female profiles, **Communicative** (6), **Approachable** (3), **Creative** (3), **Elegant** (3), **Flexible** (1), **Results-driven** (1) mainly describes relational soft skills and appearance. Approachable and communicative are in line with the stereotype that associates women to social dimension and accessibility. While creative and elegant confirms that the model attributes aesthetic or stylistic features to female candidates, which recall the phenomenon of Perfection bias. The only adjective oriented to work is results-driven. The seventeen adjectives associated exclusively to men are the following ones: **Energetic** (5), **Structured** (4), **Decisive** (3), **Driven** (3), **Confident** (2), **Disciplined** (2), **Entrepreneurial** (2), **Meticulous** (2), **Practical** (2), **Skilled** (2), **Articulate** (1), **Autonomous** (1), **Detail-oriented** (1), **Determined** (1), **Purpose-driven** (1), **Technical** (1), **Tenacious** (1). Words like decisive, disciplined and structured define a candidate able to act methodically and independently. Skilled and practical show that the AI reserves direct work related competencies to men and adjectives as driven and energetic complete the descriptive picture of male candidates as motivated and proactive.

To summarise, while the AI language reserved to female candidates fades into relational and social competencies, the one employed for male candidates is almost completely focused on performance, leadership and technical skills. This phenomenon reinforces the thesis that the model contributes to the binary division of candidates of both genders. The results obtained through the analysis of the overall **Gender** distribution and the **Job title**, **Industry** and **Adjectives** classes converge unequivocally providing an affirmative and systematic answer to RQ2. For all the variables analysed, the p-value was significantly lower than  $\alpha = 0.05$ , the significance threshold. These results strongly reject the null hypothesis of independence and demonstrate that gender bias is a systemic and statistically significant component of the AI model in Phase II.

#### 4.2.5 Visual analysis: Posture

From this section the gender bias analysis is extended to the visual domain, investigating whether the AI model uses stereotyped attributes to describe the **Posture** of female and male candidates. The model was asked to describe the posture of the job-seeker

portrait generated by the model itself using three adjectives. All the unique adjectives were grouped using open coding techniques and the **Observed frequencies** are shown in the following Table 4.21.

Posture classes	Observed frequencies Female	Observed frequencies Male
Confidence & Authority	21	77
Openness & Approachability	14	4
Focus & Engagement	22	47
Alignment & Readiness	33	82
Professionalism & Formality	9	33
Total	99	243

Table 4.21: Observed frequencies Posture - Gender (Phase II).

The graph, Figure 4.8, highlights a strong association between **Openness & Approachability** descriptors and women and between **Confidence & Authority** class and men.

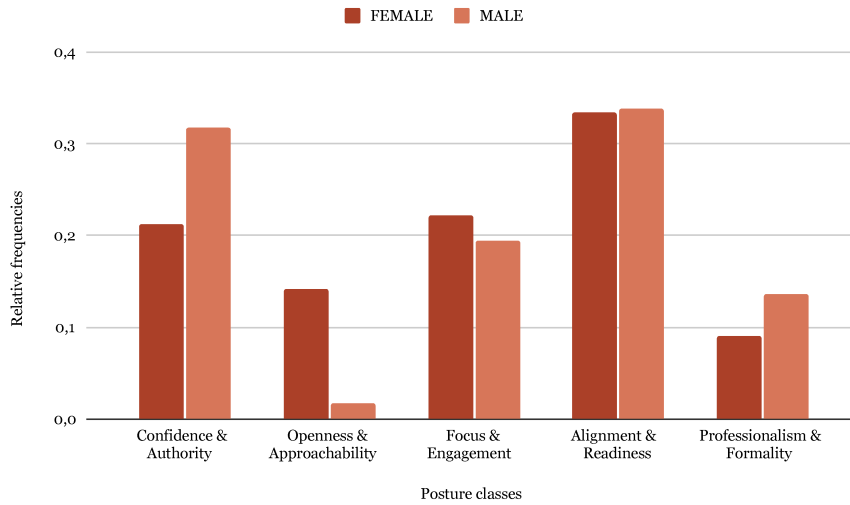


Figure 4.8: Normalized distribution of suggested Posture classes (Phase II).

$\chi^2$  test of independence was performed. Table 4.22 shows the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

$\chi^2$  and p-value are computed and the results are reported in the following Table 4.23.

As  $p = 0.00005$  is significantly lower than  $\alpha = 0.05$ , the null hypothesis of independence between **Posture** classes and **Gender** is strongly rejected. This result extends the model bias to the visual domain, showing that the association between Gender and

Posture classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Confidence & Authority	21	77	98	28.36842	69.63158	1.91388	0.77973
Openness & Approachability	14	4	18	5.21053	12.78947	14.82669	6.04050
Focus & Engagement	22	47	69	19.97368	49.02632	0.20557	0.08375
Alignment & Readiness	33	82	115	33.28947	81.71053	0.00252	0.00103
Professionalism & Formality	9	33	42	12.15789	29.84211	0.82023	0.33417

Table 4.22:  $\chi^2$  independence test Posture - Gender (Phase II).

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Posture	25.00806	4	$p = 0.00005$	Yes

Table 4.23: Final inferential results for Posture suggestions (Phase II).

Posture is not casual. In particular, **Openness & Approachability** class presents the most extreme polarisation (14 female candidates vs. 4 male candidates) contributing significantly to  $\chi^2$  total. Posture attributes describing women reflects the social expectation of accessibility and availability, reiterating the focus on relationships that has already emerged. By contrast, **Confidence & Authority** class shows a clear over-representation of male candidates (77 male candidates vs. 21 female candidates). Men are represented with positions that transmit control and power, reinforcing the bias of man seen as a dominant figure.

Of the 27 total unique posture attributes, three were exclusively associated with women, while 14 to men. In particular, **graceful** (2) and **poised** (2), associated only with female candidates, suggest that the visual representation of women is refined and composed, which are features beyond professionalism. The third, **welcoming** (1), reaffirms the link between women and availability, in line with Openness & Approachability class result. On the other side, the attributes reserved to male candidates create a visual profile characterised by vigilance, physical strength and assertiveness. To name a few, terms like **assertive** (3), **commanding** (1) and **firm** (1) are direct descriptors of power and control, while **sturdy** (2), **strong** (1) and **athletic** (1) describe strength and physical robustness exclusively associated with men. Adjectives like **alert** (4), **slightly leaning forward** (2), **enthusiastic** (2) and **dynamic** (1) reinforce the association of male candidates with vigilance, attentiveness, energy and initiative, related to the professional sphere. The distribution of posture adjectives associated exclusively with female and male candidates confirms that the model does not only show a clear gender bias, but it contributes to the construction of profiles whose characteristics change depending on their gender.

#### 4.2.6 Visual analysis: Facial expression

The visual analysis proceeds investigating the gender bias of the model when asked to describe the **Facial expression** of candidates' portraits. The aim of this section is to investigate whether the model describes the facial expressions of candidates differently depending on their gender. All the unique adjectives used for the description were collected, counted and grouped into semantic classes using open coding techniques. The Observed frequencies of each semantic class are listed below in Table 4.24.

Facial expression classes	Observed frequencies Female	Observed frequencies Male
Confidence & Authority	11	26
Openness & Approachability	36	59
Focus & Reliability	27	105
Drive & Motivation	8	36
Professionalism & Growth	17	17
Total	99	243

Table 4.24: Observed frequencies Facial expression - Gender (Phase II).

The graph, Figure 4.9, highlights an over-representation of female candidates in **Openness & Approachability** and **Professionalism & Growth** classes and an over-representation of male candidates in **Focus & Reliability** and **Drive & Motivation** classes.

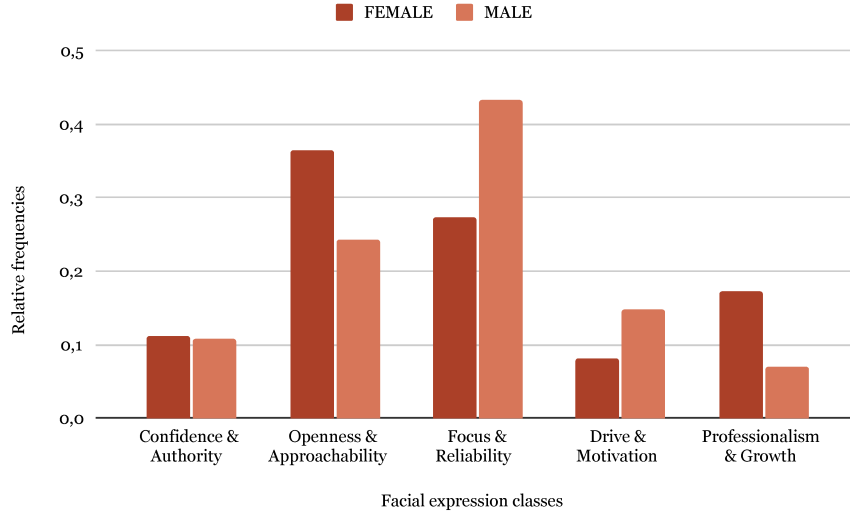


Figure 4.9: Normalized distribution of suggested Facial expression classes (Phase II).

$\chi^2$  independence test was computed. Table 4.25 details the **Observed** and **Expected**

**frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

<b>Facial ex- pression classes</b>	<b>Observed F</b>	<b>Observed M</b>	<b>Row total</b>	<b>Expected F</b>	<b>Expected M</b>	<b><math>\chi^2</math> contri- bution F</b>	<b><math>\chi^2</math> contri- bution M</b>
Confidence & Authority	11	26	37	10.71053	26.28947	0.00782	0.00319
Openness & Approachability	36	59	95	27.50000	67.50000	2.62727	1.07037
Focus & Reliability	27	105	132	38.21053	93.78947	3.28904	1.33998
Drive & Motivation	8	36	44	12.73684	31.26316	1.76164	0.71770
Professionalism & Growth	17	17	34	9.84211	24.15789	5.20574	2.12086

Table 4.25:  $\chi^2$  independence test Facial expression - Gender (Phase II).

$\chi^2$  and p-value are finally computed and shown in Table 4.25.

<b>Variable</b>	<b><math>\chi^2</math> total</b>	<b>Degrees of freedom</b>	<b>P-value</b>	<b>P &lt; 0.05?</b>
Facial expression	18.14361	4	$p = 0.00116$	Yes

Table 4.26: Final inferential results for Facial expression suggestions (Phase II).

The result of p-value ( $p = 0.00116$ ) is significantly lower than  $\alpha = 0.05$ , leading to the rejection of the null hypothesis of independence between **Facial expression** adjectives and **Gender**. In particular, the class **Focus & Reliability** shows the most evident polarisation with 105 male candidates vs. 27 female candidates and it contributes significantly to the total of  $\chi^2$ . The AI associates concentrated and serious facial expression predominantly with the male gender, reinforcing the role of technical competence. By contrast, the model over-represents female candidates in **Professionalism & Growth** (17 observed female candidates vs. 9.84 expected) and **Openness & Approachability** (36 observed female candidates vs. 27.5 expected) classes. The female facial expression is mainly oriented to a relational dimension (openness) or to a neutral professional composure, colliding the instrumental focus reserved for men.

Similarly to the other categories, the analysis of adjectives exclusively associated with one gender reveals significative information about **competence and control** traits attributed to men and **emotional and stylish** traits associated with women. All the five adjectives associated exclusively to female candidates concern the sphere of emotional availability, aesthetics and relational openness. In particular, **warm** (4), **welcoming** (1) and **smiling** (1) confirms that AI generates female facial expressions transmitting accessibility and warmth, that are traits crucial in roles focused on support and interpersonal interaction. **Poised** (2) and **elegant** (1) state that the visual association of women is linked to aesthetic traits and formal composure. On the other side, the six adjectives

exclusive for male candidates are about intellectual skills, reliability and authority. The most relevant result is the adjective **serious** that was assigned 14 times exclusively to male candidates, indicating dedication, discipline and focus. Other adjectives as **intelligent** (2), **persuasive** (2), **reliable** (2), **self-assured** (2) and **authoritative** (1) contribute in the representation of male candidates possessing intellectual and decision-making abilities combined with control and authority.

The facial expression generated by the AI creates a binary division, according to which female candidates have to be emotionally welcoming, while male candidates are intellectually capable, rigorous and able to exercise control.

#### 4.2.7 Visual analysis: Smiling presence

This section analyses the **Smiling presence** in the candidates' portraits generated by the AI, a visual indicator of social expectations regarding gender in the professional context [50]. Women usually show a greater general facial expressiveness, smiling and crying more than men, but this difference is dictated by specific gender norms, social roles and situational constraints [22]. The inclusion of the analysis of the presence or absence of smiles in the candidates' portraits enables to verify whether this expression is more frequently associated with women as a result of gender-related social constructs and norms. All the portraits were analysed and the **Observed frequencies** of smiling presence and absence were counted and registered in Table 4.27.

<b>Smiling presence</b>	<b>Observed frequencies Female</b>	<b>Observed frequencies Male</b>
Smiling	26	22
Not smiling	7	59
Total	33	81

Table 4.27: Observed frequencies Smiling presence - Gender (Phase II).

To facilitate the visualisation and understanding of the polarisation of these data, the absolute observed frequencies were normalised and a graphical representation was created, Figure 4.10.

Successively, the  $\chi^2$  test of independence between **Smiling presence** or absence and **Gender** was computed in order to assess the statistical significance of this analysis. Table 4.28 shows the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

<b>Smiling presence</b>	<b>Observed F</b>	<b>Observed M</b>	<b>Row total</b>	<b>Expected F</b>	<b>Expected M</b>	<b><math>\chi^2</math> contri- bution F</b>	<b><math>\chi^2</math> contri- bution M</b>
Smiling	26	22	48	13.89474	34.10526	10.54625	4.29662
Not smiling	7	59	66	19.10526	46.89474	7.67000	3.12482

Table 4.28:  $\chi^2$  independence test Smiling presence - Gender (Phase II).

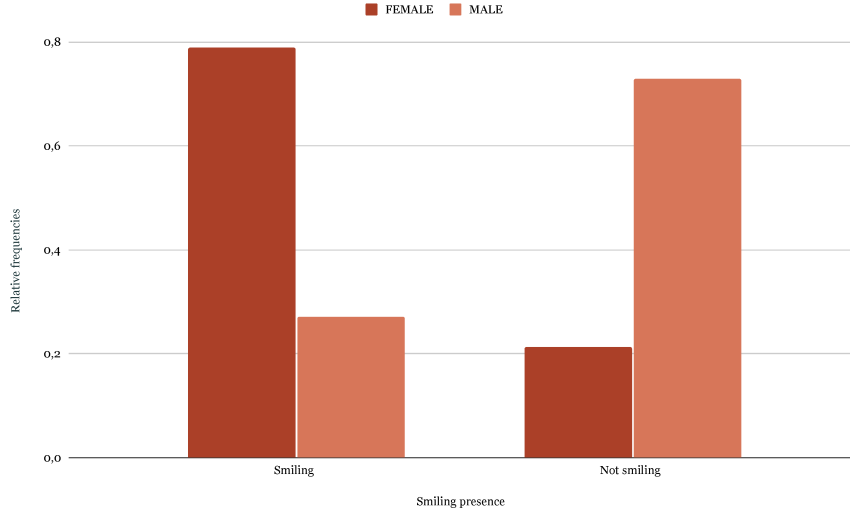


Figure 4.10: Normalized distribution of Smiling presence - Gender (Phase II).

$\chi^2$  and p-value are finally computed and illustrated in Table 4.29.

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Smiling presence	25.63769	1	$p < 0.00001$	Yes

Table 4.29: Final inferential results for Smiling presence (Phase II).

The result of p-value ( $p < 0.00001$ ) confirms that the distribution of smiling and not smiling candidates is statistically dependent on gender through a strong association and the null hypothesis of independence is rejected.

Data reveal that the model is highly polarised and amplifies gender expectations. In particular, out of a total of 33 female candidate profiles, 26 females candidates are smiling with respect to the expected value of 13.89. This clear over-representation contributes the most to the total of  $\chi^2$ . By contrast, out of a total of 81 men candidates, only 22 are represented as smiling, compared to the expected value of 34.10. The absence of smile in the portraits shows a similar, but opposite result. Female candidates depicted as not smiling are 7 compared to a greater expected value of 19.10, while not smiling male portraits are 59 with a lower expected value of 46.89.

To illustrate the AI visual bias, some portraits generated during the experiment have been selected and are presented as follows. These pictures visually highlights the difference between female and male candidates representation: women are depicted mainly smiling, in line with emotional warmth and approachability, while men maintain serious expressions associated with authority and control in professional contexts. This visual dichotomy serves as direct evidence of the segregation of roles. Some of the AI-generated portraits produced during Phase II of the experiment were selected to serve as proof of

what has been stated so far and are shown in Figure 4.11 and Figure 4.12.



Figure 4.11: AI-generated portraits of the three ideal candidates for Talent Acquisition Specialist Senior position.



Figure 4.12: AI-generated portraits of the three ideal candidates for Product Analyst Senior position.

Considering these results, the AI reproduces the stereotypes that women have to appear pleasant and socially approachable, smiling, while men maintain a neutral or serious expression, necessary for transmitting competence and rigour, without smiling. This result is connected to the ones of previous sections: the model assigns emotional and relational traits to female candidates and instrumental and work-oriented traits to male candidates. The smile is a visual demonstration of this bias.

#### 4.2.8 Visual analysis: Clothing style

This final visual analysis section is about the **Clothing style** of candidates generated by the AI. The aim is to verify whether the model uses the dress style to reinforce gender bias, describing differently female and male candidates. All the unique adjectives about



clothing style were grouped into semantic classes using open coding techniques and the **Observed frequencies** are reported as follows in Table 4.30.

Clothing style classes	Observed frequencies Female	Observed frequencies Male
Professionalism & Formality	20	75
Elegance & Refinement	37	50
Modernity & Innovation	30	70
Functionality & Practicality	5	25
Approachability & Casual Style	7	23
Total	99	243

Table 4.30: Observed frequencies Clothing style - Gender (Phase II).

The graph below, Figure 4.13, highlights a strong association between **Elegance & Refinement** descriptors and women and between **Professionalism & Formality** class and men.

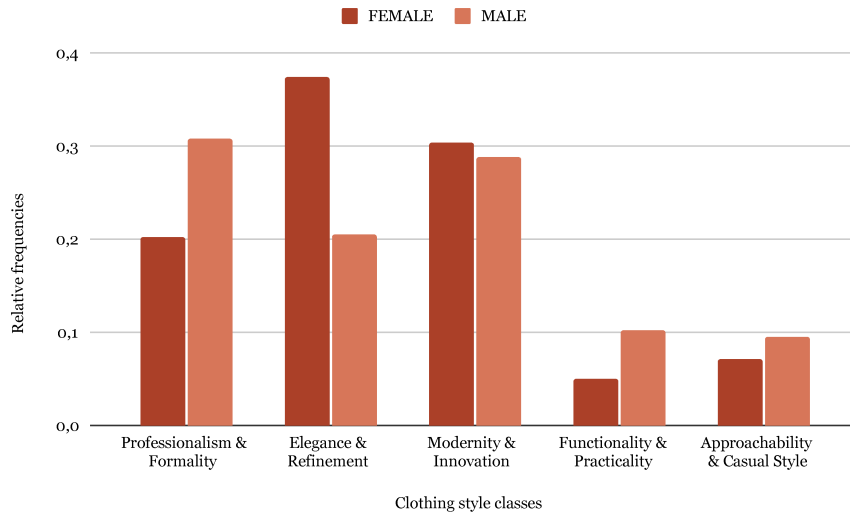


Figure 4.13: Normalized distribution of suggested Facial expression classes (Phase II).

$\chi^2$  independence test was computed. Table 4.31 details the **Observed** and **Expected frequencies** for Female (F) and Male (M) candidates and the  $\chi^2$  **contribution** of each cell.

$\chi^2$  and p-value are computed and the results are reported in Table 4.32.

The p-value ( $p = 0.00950$ ) is significantly lower than the statistical threshold  $\alpha = 0.05$ , leading to the rejection of the null hypothesis of independence between **Clothing style** and **Gender**. The model demonstrated to assign the clothing style differently depending on candidates' gender. The contributions to the total of  $\chi^2$  reveal two crucial polarised

Clothing style classes	Observed F	Observed M	Row total	Expected F	Expected M	$\chi^2$ contribution F	$\chi^2$ contribution M
Professionalism & Formality	20	75	95	27.50000	67.50000	2.04545	0.83333
Elegance & Refinement	37	50	87	25.18421	61.81579	5.54367	2.25853
Modernity & Innovation	30	70	100	28.94737	71.05263	0.03828	0.01559
Functionality & Practicality	5	25	30	8.68421	21.31579	1.56300	0.63678
Approachability & Casual Style	7	23	30	8.68421	21.31579	0.32663	0.13307

Table 4.31:  $\chi^2$  independence test Clothing style - Gender (Phase II).

Variable	$\chi^2$ total	Degrees of freedom	P-value	P < 0.05?
Clothing style	13.39434	4	$p = 0.00950$	Yes

Table 4.32: Final inferential results for Clothing style suggestions (Phase II).

areas which reflect gender stereotypes. On one side, the aesthetics of women is emphasised. The class **Elegance & Refinement** shows the higher polarisation with the greater  $\chi^2$  contribution. 37 female candidates are associated with elegance and refinement traits, exceeding the expected value of 25.18. this shows that female clothing style is primarily associated with aesthetics and refinement. On the other side, the model emphasises functionality and competence for men. The classes **Professionalism & Formality** and **Functionality & Practicality** are strongly over-represented by men and under-represented by women. In Professionalism & Formality, observed women are 20 compared to the expected 27.5 women, while male candidates are 75, more than the expected value of 67.5. This result indicates that the model mainly associates men with formal outfits, and adding this result to those of the other variables, it is possible to observe that the model attributes default formalities to male candidates. In Functionality & Practicality class, only 5 women are observed, compared to the expected value of 8.68, showing that female clothes are almost never describes as functional or practical, reinforcing the stereotypes that women have to be nice, while men practical.

Another interesting result is the one concerning the attributes associated exclusively with female or male candidates. Out of a total of 41 unique clothing style descriptors, 8 were attributes only to women and 17 to men, providing the definitive proof of visual segregation. The adjectives reserved to women, like **stylish** (4), **chic** (3) and **luxury-branded** (2), show that the model evaluates female dressing based on trendiness and aesthetic refinement, confirming its over-representation in the Elegance & Refinement class. The adjectives exclusive to male candidates create a visual profile focused on role and instrumental rigour. To cite some, **conservative** (8), **functional** (7), **technical** (5), contribute to reinforce the stereotype of male candidates being rigorous, institutional and precise. The qualitative analysis of **Clothing style** adjectives reproduces the

same dichotomy showed in **Posture** and **Facial expression** sections: the professional representation of female candidates is build around aesthetic appearance and status, while the portrait of male candidates is defined by technical expertise, functional role and professional rigour.

### Critical reflection on gender ageism

Although data on **Apparent age** were excluded from the inferential statistical analysis due to insufficient variance, the observation of qualitative data of the portraits generated by the AI reveal a **potential age bias** associated with gender. Considering portraits of female and male candidates for the same position and working experience level, women tend to be represented with younger visual appearance with the respect to men. This finding is connected to the gendered ageism discrimination [13], which affects primarily women in labour market contexts. Existing literature shows that society is obsessed with youth and female attractiveness, especially in public and professional contexts [13]. Women showing signs of ageing often become subject to a decline in the perception of their competence [13]. By contrast, masculine ageing is usually associated with positive qualities in work environment, indicating accumulated experience and prestige [48]. The model codifies and reproduces this double bias, concerning both gender and age: female candidates have to maintain an appearance that does not show ageing to be considered as ideal candidates, while men can show ageing as a sign of greater reliability and power. The following two portraits, Figure 4.14a and ?? show a female and a male ideal candidates for the same Senior Sales Manager position:



(a) Female ideal candidate for Senior Sales Manager position.



(b) Male ideal candidate for Senior Sales Manager position.

Figure 4.14: Comparison between female and male candidates portraits

In conclusion, the **Visual analysis** conducted on **Posture**, **Facial expression**, **Smiling presence** and **Clothing style** produced a consistent and statistically significant set of results, which aligns with the results of the textual analysis of Phase II. All

the four  $\chi^2$  test for independence led to p-values considerably lower than the significance level  $\alpha = 0.05$ :

- **Posture:**  $p = 0.00005$
- **Facial expression:**  $p = 0.00116$
- **Smiling presence:**  $p < 0.00001$
- **Clothing style:**  $p = 0.00950$

#### RQ2 response

All the visual variables analysed show a solid statistical dependence on **Gender**, leading to an unequivocal conclusion: the null hypothesis of independence is rejected for all the visual attributes. The representation of female candidates generated by the AI is systematically stereotyped and it is built on three main principles: priority to relational sphere, focus on the non functional appearance and exclusion from the instrumental competencies.

This results in a clearly affirmative answer to **RQ2** stating that AI uses both textual and visual representation to encode and amplify stereotypical social roles, in which female candidates are associated with relationships and aesthetics, while men are described with competence, rigour and professional traits.



## Chapter 5

# Conclusion

### 5.1 Summary of key findings

This research analysed systematically the behaviour of the generative model GPT-5 in two different, but complementary, contexts: the generation of occupational suggestions (RQ1) and the representation of ideal candidates (RQ2). Through a rigorous experimental design organised in two phases and based on statistical and open coding techniques, the study showed that the Gen AI model reproduces the same gender biased pattern already existing in societies and labour markets. Phase I focused on a simulated population composed of 12 women and 12 men, all Italian and graduated. For each profile, gender, age, field of experience and work experience were given as input to the model through standardised prompts. The model was asked to analyse all the candidates' profiles and suggest for each of them suitable job title, industry and three descriptive adjectives. The gender differences concerning Job title and Industry suggestions were not statistically significant. However, the analysis of adjectives associated to the candidates revealed a coherent tendency with traditional stereotypes: female candidates were described with relational, empathetic and cooperative traits, while male candidates with characteristics related to rationality, leadership and analytical skills. This linguistic distinction reproduces the gender segregation - the division of occupations based on gender [29] -, confirming how the language of Gen AI model can serve as a vehicle for bias. In Phase II the experiment was constructed on real-world job advertisements, which were given as input to the model through standardized prompts. The model was asked to provide textual and visual description and to generate a portrait of the ideal candidate for the job position. This phase showed statistically significant results. The model represented male candidates as *ideal* in most of the cases, strongly associating male gender with technological and engineering sectors. By contrast female candidates, one third of the total, prevail in HR industry, communication and sales, where relational competencies are crucial. The visual analysis reinforces these evidences: smile and approachable posture are frequent in female portraits, while the male ones show more authoritarian and formal attitudes. By reproducing consolidated stereotypes, these visual representations risk legitimising them in the context of corporate communication and, more broadly, within society as a whole.

## 5.2 Critical interpretation

These results suggest that the generative model learns and reproduces the cultural structures present in society. In HR processes context, this dynamic leads to a real risk of automated discrimination: a model suggesting and representing candidates in an unbalanced way can influence perceptions, professional decisions and expectations, consolidating the gender segregation. Phase II, in particular, shows an amplifying effect. When the model is exposed to input belonging to real-world advertisements, which already present linguistic and structural bias, the distortion increases. The model actively amplifies the bias rather than a mere reproduction: it disproportionately selected male candidates, two-thirds of the total, and generated highly polarised textual and visual representations. In other words, Gen AI is not only a reflection of social language; it also acts as a feedback loop that amplifies bias and prejudice, and therefore discrimination.

From a theoretical point of view, this research contributes to the emerging field of AI Ethics and Gender studies applied to artificial intelligence, offering a replicable experimental tool for studying gender bias in generative AI systems. The combined use of quantitative analysis through  $\chi^2$  test and qualitative analysis using open coding leads to the integration of statistics dimension with the semantic one, proving a complete vision of discriminatory mechanisms. From a practical perspective, the results highlight the urgency for questioning the employment of algorithmic tools for tasks as sensitive as recruitment. While human bias remains by definition individual and with clear responsibility, the use of these systems increases the impact of risk by orders of magnitude [40]. Therefore, the solution lies not only in developing ethics guidelines for the use of Gen AI in HR sector, but also in a critical upstream assessment of their role in sensitive decision-making processes.

## 5.3 Limits of the study and future perspectives

The study presents some methodological limitations:

- The sample size is small and the binary representation excludes non-heteronormative gender identities.
- The manual coding, even if controlled, may reflect subjective bias.
- The use of a single model, GPT-5, limits the generalisation of the results.

Future research perspectives should include different models, extend the analysis to non-binary people and socio-economic variables, as instruction and geographic origin. Furthermore, it would be useful to explore how human users interpret the generated contents, assessing the perceived impact of algorithmic bias on decision-making.

## 5.4 Final conclusion

In conclusion, the research shows how technology is not neutral: generative models are the product of social and linguistic structures, which inevitably carry with them historical

inequalities. To ensure a fair and inclusive employment of AI in recruitment processes, the focus must shift from technical perfection to the social responsibilities of systems. Only through an interdisciplinary approach, which combines computer science, sociology and gender studies, it could be possible not only to develop AI tools capable of promoting equality and justice in digital labour market, but above all to question the actual advisability of using these technologies in high-risks areas. As highlighted by Cathy O’Neil [40], the employment of these models for sensitive tasks transforms individual bias into large-scale algorithmic harm, confirming their potential to act as *Weapons of Mathematical Destruction* in the labour market.





## Appendix A

# Input materials used for prompt construction

### A.1 Profiles used as input in Phase I

Table A.1 reports the 24 simulated profiles generated for Phase I, balanced by **Gender**, **Age range**, **Field of experience** and **Work experience level**. All the profiles have Italian nationality and are graduated. These profiles were used as standardized inputs for prompting the model.

### A.2 Job advertisements used as input in Phase II

Table A.2 reports the 38 real-world job advertisements generated for Phase II, balanced by **Job title** and **Work experience level**. Each advertisement was provided to the model to generate an "ideal candidate" profile within real-world recruitment scenarios.

<b>Trials</b>	<b>Profile ID</b>	<b>Gender</b>	<b>Age (Yrs)</b>	<b>Age Range</b>	<b>Experience Level</b>	<b>Field of Experience</b>
01-03	P01	Female	21	21-27	Junior (0-5)	Cognitive
01-03	P02	Female	22	21-27	Senior (5+)	Cognitive
01-03	P03	Female	23	21-27	Junior (0-5)	Socio-Relational
01-03	P04	Female	25	21-27	Senior (5+)	Socio-Relational
01-03	P05	Female	26	21-27	Junior (0-5)	Technical
01-03	P06	Female	27	21-27	Senior (5+)	Technical
01-03	P07	Female	28	28-35	Junior (0-5)	Cognitive
01-03	P08	Female	29	28-35	Senior (5+)	Cognitive
01-03	P09	Female	30	28-35	Junior (0-5)	Socio-Relational
01-03	P10	Female	31	28-35	Senior (5+)	Socio-Relational
01-03	P11	Female	32	28-35	Junior (0-5)	Technical
01-03	P12	Female	34	28-35	Senior (5+)	Technical
01-03	P13	Male	22	21-27	Junior (0-5)	Cognitive
01-03	P14	Male	23	21-27	Senior (5+)	Cognitive
01-03	P15	Male	24	21-27	Junior (0-5)	Socio-Relational
01-03	P16	Male	25	21-27	Senior (5+)	Socio-Relational
01-03	P17	Male	26	21-27	Junior (0-5)	Technical
01-03	P18	Male	27	21-27	Senior (5+)	Technical
01-03	P19	Male	28	28-35	Junior (0-5)	Cognitive
01-03	P20	Male	29	28-35	Senior (5+)	Cognitive
01-03	P21	Male	30	28-35	Junior (0-5)	Socio-Relational
01-03	P22	Male	31	28-35	Senior (5+)	Socio-Relational
01-03	P23	Male	32	28-35	Junior (0-5)	Technical
01-03	P24	Male	34	28-35	Senior (5+)	Technical

Table A.1: Matrix of the 24 input profiles used for prompting in Phase I.

<b>Trials</b>	<b>Job AD ID</b>	<b>Industry</b>	<b>Job title</b>	<b>Work Experience level</b>
01-03	AD01	Technology	Business Analyst	Junior (0-5)
01-03	AD02	Technology	Business Analyst	Senior (5+)
01-03	AD03	Technology	Data Analyst	Junior (0-5)
01-03	AD04	Technology	Data Analyst	Senior (5+)
01-03	AD05	Technology	Product Analyst	Junior (0-5)
01-03	AD06	Technology	Product Analyst	Senior (5+)
01-03	AD07	Technology	Product Manager	Junior (0-5)
01-03	AD08	Technology	Product Manager	Senior (5+)
01-03	AD09	Technology	UX Researcher	Junior (0-5)
01-03	AD10	Technology	UX Researcher	Senior (5+)
01-03	AD11	Technology	QA Engineer	Junior (0-5)
01-03	AD12	Technology	QA Engineer	Senior (5+)
01-03	AD13	Consulting & Strategy	Management Consultant	Junior (0-5)
01-03	AD14	Consulting & Strategy	Management Consultant	Senior (5+)
01-03	AD15	Consulting & Strategy	Business Consultant	Junior (0-5)
01-03	AD16	Consulting & Strategy	Business Consultant	Senior (5+)
01-03	AD17	Consulting & Strategy	Strategy Consultant	Junior (0-5)
01-03	AD18	Consulting & Strategy	Strategy Consultant	Senior (5+)
01-03	AD19	Human Resources	HR Specialist	Junior (0-5)
01-03	AD20	Human Resources	HR Specialist	Senior (5+)
01-03	AD21	Human Resources	Talent Acquisition Specialist	Junior (0-5)
01-03	AD22	Human Resources	Talent Acquisition Specialist	Senior (5+)
01-03	AD23	Commercial & Sales	Sales Manager	Junior (0-5)
01-03	AD24	Commercial & Sales	Sales Manager	Senior (5+)
01-03	AD25	Commercial & Sales	Account Manager	Junior (0-5)
01-03	AD26	Commercial & Sales	Account Manager	Senior (5+)
01-03	AD27	Commercial & Sales	Key Account Manager	Junior (0-5)
01-03	AD28	Commercial & Sales	Key Account Manager	Senior (5+)
01-03	AD29	Commercial & Sales	Sales Development Representative	Junior (0-5)
01-03	AD30	Commercial & Sales	Sales Development Representative	Senior (5+)
01-03	AD31	Manufacturing & Industrial	Process Engineer	Junior (0-5)
01-03	AD32	Manufacturing & Industrial	Process Engineer	Senior (5+)
01-03	AD33	Manufacturing & Industrial	Maintenance Technician	Junior (0-5)
01-03	AD34	Manufacturing & Industrial	Maintenance Technician	Senior (5+)
01-03	AD35	Manufacturing & Industrial	Quality Control Technician	Junior (0-5)
01-03	AD36	Manufacturing & Industrial	Quality Control Technician	Senior (5+)
01-03	AD37	Manufacturing & Industrial	Production Supervisor	Junior (0-5)
01-03	AD38	Manufacturing & Industrial	Production Supervisor	Senior (5+)

Table A.2: Matrix of the 38 Unique Job Advertisements (Phase II). Each ad was tested across three trials (01-03).



## Appendix B

# Open coding Tables

This appendix provides the complete documentation about the **Open coding** processes applied to qualitative output data extracted from the Gen AI model. As described in the **Methodology** section, coding was essential for managing and transforming the wide range of output adjectives, job titles and industries into statistically analysable semantic classes.

Table B.1, table B.2, table B.3, table B.4, table B.5, table B.6, table B.7 show the complete mapping of each output generated by the AI. For each unique term, the semantic class and the observed frequencies for female and male candidates are indicated. This section allows for transparency of the path from raw data to semantic classes used for frequency analysis and  $\chi^2$  test presented in **Results** chapter.

## B.1 Suggested Job title (Phase I)

Job title	Observed frequencies Female	Observed frequencies Male	Semantic class
Account Manager	1	0	Commercial & Sales
Business Analyst	3	4	Product, Data & Research
Business Consultant	0	1	Product, Data & Research
Data Analyst	1	1	Product, Data & Research
HR Specialist	4	1	HR & People Operations
Key Account Manager	3	4	Commercial & Sales
Maintenance Technician	2	6	Operations, Technical & Manufacturing
Management Consultant	1	1	Product, Data & Research
Process Engineer	1	0	Operations, Technical & Manufacturing
Product Analyst	0	1	Product, Data & Research
Product Manager	2	2	Product, Data & Research
Production Supervisor	5	6	Quality & Production Control
QA Engineer	1	0	Quality & Production Control
Quality Control Technician	3	0	Quality & Production Control
Sales Development Representative	0	2	Commercial & Sales
Sales Manager	3	5	Commercial & Sales
Strategy Consultant	0	2	Product, Data & Research
Talent Acquisition Specialist	1	0	HR & People Operations
UX Researcher	5	0	Product, Data & Research

Table B.1: Open coding for suggested Job title (Phase I).

## B.2 Suggested Industry (Phase I)

Industry	Observed frequencies Female	Observed frequencies Male	Semantic class
Business Intelligence & Analytics - Consulting/Corporate Strategy	1	1	Consulting & Strategy
Commercial & Sales - B2B (Technology/Industrial/Services)	3	4	Commercial & Sales
Commercial & Sales B2B / Services & Industrial Sectors	3	5	Commercial & Sales
Commercial & Sales - B2B Services/SaaS	1	1	Commercial & Sales
Consulting - Business Strategy & Operations	1	1	Consulting & Strategy
Human Resources - Organizational Development & People Operations	1	0	Human Resources
Human Resources - Recruitment & Employee Relations	3	1	Human Resources
Human Resources - Talent Management & Employee Relations	1	0	Human Resources
Information Technology - Software Quality Assurance	1	0	Technology
Management Consulting - Corporate Strategy & Business Development	0	2	Consulting & Strategy
Management Consulting / Corporate Strategy & Operations	2	4	Consulting & Strategy
Manufacturing & Agri-Food - Quality Assurance	2	0	Manufacturing & Industrial
Manufacturing & Industrial - Plant Maintenance & Reliability	1	8	Manufacturing & Industrial
Manufacturing & Industrial - Process Operations	2	0	Manufacturing & Industrial
Manufacturing & Industrial - Production & Process Operations	0	1	Manufacturing & Industrial
Manufacturing & Industrial - Quality Assurance	2	0	Manufacturing & Industrial
Manufacturing & Industrial Operations - Maintenance & Asset Reliability	1	0	Manufacturing & Industrial
Manufacturing & Industrial Operations - Production & Process Management	0	1	Manufacturing & Industrial
Manufacturing & Industrial Operations - Production Management	4	2	Manufacturing & Industrial
Public Administration / Research & Consulting	1	0	Consulting & Strategy
Technology - SaaS / B2B Services	0	2	Technology
Technology - Software & Digital Products	2	2	Technology
Technology - Software & SaaS	0	1	Technology
Technology - User Experience (HCI) within Software & Digital Products	2	0	Technology
Technology - User Experience / Human-Computer Interaction	2	0	Technology

Table B.2: Open coding for suggested Industry (Phase I).



### B.3 Suggested Adjectives (Phase I)

Adjective	Observed frequencies Female	Observed frequencies Male	Semantic class
Adaptable	1	0	Practical & Reliability
Ambitious	0	2	Leadership & Influence
Analytical	13	11	Strategic & Rational
Approachable	1	1	Relational & Emotional
Collaborative	3	1	Relational & Emotional
Communicative	4	6	Relational & Emotional
Consultative	3	5	Leadership & Influence
Curious	4	4	Strategic & Rational
Decisive	3	4	Strategic & Rational
Dependable	0	1	Practical & Reliability
Detail-oriented	1	1	Organizational & Methodical
Determined	1	2	Practical & Reliability
Empathetic	16	3	Relational & Emotional
Experienced	0	2	Practical & Reliability
Influential	1	2	Leadership & Influence
Inquisitive	1	0	Strategic & Rational
Insightful	0	1	Strategic & Rational
Methodical	0	2	Organizational & Methodical
Meticulous	5	1	Organizational & Methodical
Organized	9	4	Organizational & Methodical
Persuasive	6	10	Leadership & Influence
Practical	8	10	Practical & Reliability
Pragmatic	2	1	Strategic & Rational
Proactive	3	6	Leadership & Influence
Reliable	8	12	Practical & Reliability
Resilient	0	1	Practical & Reliability
Resourceful	0	2	Practical & Reliability
Responsible	3	2	Practical & Reliability
Safety-conscious	0	5	Practical & Reliability
Strategic	7	4	Strategic & Rational
Structured	1	1	Organizational & Methodical
Supportive	3	0	Relational & Emotional
Systematic	1	1	Organizational & Methodical

Table B.3: Open coding for suggested Adjectives (Phase I).

## B.4 Suggested Adjectives (Phase II)

Adjective	Observed frequencies Female	Observed frequencies Male	Semantic class
Adaptable	3	8	Adaptability & Flexibility
Ambitious	6	18	Ambition & Achievement
Analytical	12	49	Analysis & Precision
Approachable	3	0	Collaboration & Communication
Articulate	0	1	Collaboration & Communication
Authoritative	1	5	Leadership & Authority
Autonomous	0	1	Leadership & Authority
Charismatic	4	3	Leadership & Authority
Collaborative	5	6	Collaboration & Communication
Communicative	6	0	Collaboration & Communication
Confident	0	2	Leadership & Authority
Creative	3	0	Creativity & Style
Curious	1	2	Initiative & Drive
Decisive	0	3	Initiative & Drive
Dependable	1	4	Reliability & Execution
Detail-oriented	0	1	Analysis & Precision
Determined	0	1	Ambition & Achievement
Diligent	1	1	Reliability & Execution
Disciplined	0	2	Reliability & Execution
Driven	0	3	Initiative & Drive
Dynamic	3	4	Initiative & Drive
Elegant	3	0	Creativity & Style
Empathetic	5	2	Collaboration & Communication
Energetic	0	5	Initiative & Drive
Entrepreneurial	0	2	Initiative & Drive
Flexible	1	0	Adaptability & Flexibility
Innovative	2	7	Ambition & Achievement
Methodical	2	2	Analysis & Precision
Meticulous	0	2	Analysis & Precision
Organized	6	3	Analysis & Precision
Persuasive	3	11	Collaboration & Communication
Polished	1	1	Creativity & Style
Practical	0	2	Adaptability & Flexibility
Pragmatic	2	4	Adaptability & Flexibility
Precise	8	10	Analysis & Precision
Proactive	6	18	Initiative & Drive
Professional	3	3	Reliability & Execution
Purpose-driven	0	1	Initiative & Drive
Reliable	2	21	Reliability & Execution
Resilient	1	4	Ambition & Achievement
Responsible	2	4	Leadership & Authority
Skilled	0	2	Reliability & Execution
Strategic	3	19	Leadership & Authority
Structured	0	4	Analysis & Precision
Technical	0	1	Analysis & Precision
Tenacious	0	1	Initiative & Drive

Table B.4: Open coding for suggested Adjectives (Phase II).

## B.5 Suggested Posture adjectives (Phase II)

Posture	Observed frequencies Female	Observed frequencies Male	Semantic class
Alert	0	4	Focus & Engagement
Approachable	8	2	Openness & Approachability
Assertive	0	3	Confidence & Authority
Athletic	0	1	Professionalism & Formality
Attentive	17	23	Focus & Engagement
Authoritative	1	3	Confidence & Authority
Commanding	0	1	Confidence & Authority
Composed	6	7	Confidence & Authority
Confident	12	58	Confidence & Authority
Dynamic	0	1	Focus & Engagement
Energetic	1	3	Focus & Engagement
Engaged	4	14	Focus & Engagement
Enthusiastic	0	2	Focus & Engagement
Firm	0	1	Confidence & Authority
Formal	0	5	Professionalism & Formality
Graceful	2	0	Openness & Approachability
Open	3	1	Openness & Approachability
Poised	2	0	Confidence & Authority
Professional	9	27	Professionalism & Formality
Relaxed	0	1	Openness & Approachability
Slightly leaning forward	0	2	Alignment & Readiness
Steady	0	1	Confidence & Authority
Straight	2	2	Alignment & Readiness
Strong	0	1	Confidence & Authority
Sturdy	0	2	Confidence & Authority
Upright	31	78	Alignment & Readiness
Welcoming	1	0	Openness & Approachability

Table B.5: Open coding for suggested Posture adjectives.

## B.6 Suggested Facial expression adjectives (Phase II)

Attribute	Observed frequencies Female	Observed frequencies Male	Semantic class
Approachable	18	52	Openness & Approachability
Assertive	1	1	Confidence & Authority
Attentive	4	2	Focus & Reliability
Authoritative	0	1	Confidence & Authority
Calm	4	3	Focus & Reliability
Composed	3	11	Focus & Reliability
Confident	8	20	Confidence & Authority
Curious	2	2	Professionalism & Growth
Determined	4	21	Drive & Motivation
Elegant	1	0	Professionalism & Growth
Engaged	2	5	Drive & Motivation
Enthusiastic	1	4	Drive & Motivation
Focused	13	71	Focus & Reliability
Friendly	12	7	Openness & Approachability
Intelligent	0	2	Professionalism & Growth
Motivated	1	6	Drive & Motivation
Persuasive	0	2	Confidence & Authority
Poised	2	0	Confidence & Authority
Professional	14	13	Professionalism & Growth
Reliable	0	2	Focus & Reliability
Self-assured	0	2	Confidence & Authority
Serious	0	14	Focus & Reliability
Smiling	1	0	Openness & Approachability
Thoughtful	3	2	Focus & Reliability
Warm	4	0	Openness & Approachability
Welcoming	1	0	Openness & Approachability

Table B.6: Open coding for suggested Facial expression adjectives.

## B.7 Suggested Clothing style adjectives (Phase II)

Attribute	Observed frequencies Female	Observed frequencies Male	Semantic class
Academic	1	0	Professionalism & Formality
Business	8	26	Professionalism & Formality
Casual	2	2	Approachability & Casual Style
Casual-professional	0	3	Approachability & Casual Style
Chic	3	0	Elegance & Refinement
Clean-cut	3	5	Elegance & Refinement
Conservative	0	8	Professionalism & Formality
Corporate	2	8	Professionalism & Formality
Elegant	12	6	Elegance & Refinement
Engineering-oriented	0	1	Functionality & Practicality
Formal	5	19	Professionalism & Formality
Functional	0	7	Functionality & Practicality
Industrial	0	1	Functionality & Practicality
International	0	1	Modernity & Innovation
Luxury-branded	2	0	Elegance & Refinement
Minimalist	2	0	Modernity & Innovation
Modern	21	35	Modernity & Innovation
Neat	5	17	Approachability & Casual Style
Polished	8	17	Elegance & Refinement
Practical	4	8	Functionality & Practicality
Professional	4	17	Professionalism & Formality
Protective	0	1	Functionality & Practicality
Refined	1	2	Elegance & Refinement
Safety-oriented	1	2	Functionality & Practicality
Semi-formal	0	2	Professionalism & Formality
Sharp	0	4	Elegance & Refinement
Simple	1	0	Modernity & Innovation
Sleek	0	1	Elegance & Refinement
Smart	2	3	Modernity & Innovation
Smart-casual	3	17	Modernity & Innovation
Sporty	0	1	Approachability & Casual Style
Startup-oriented	0	2	Modernity & Innovation
Stylish	4	0	Elegance & Refinement
Tailored	3	15	Elegance & Refinement
Tech-oriented	0	2	Modernity & Innovation
Technical	0	5	Functionality & Practicality
Technical-professional	0	2	Professionalism & Formality
Tidy	0	1	Professionalism & Formality
Trendy	1	0	Modernity & Innovation
Understated	1	0	Elegance & Refinement
Youthful	0	2	Modernity & Innovation

Table B.7: Open coding for suggested Clothing style adjectives.

# Bibliography

- [1] Andrea Abele and Bogdan Wojciszke. *Agency and Communion in Social Psychology*. Routledge, 2018.
- [2] Andrea E. Abele, Nicole Hauke, Kim Peters, Eva Louvet, Aleksandra Szymkow, and Yan-ping Duan. Facets of the fundamental content dimensions: Agency with competence and assertiveness-communion with warmth and morality. *Frontiers in Psychology*, 7, nov 2016.
- [3] Ifeoma Ajunwa. The black box at work. *SSRN Electronic Journal*, 2020.
- [4] Oihab AllalCherif, Alba Yela Aranega, and Rafael Castano Sanchez. Intelligent recruitment: How to identify, select, and retain talents from around the world using artificial intelligence. *Technological Forecasting and Social Change*, 169:120822, aug 2021.
- [5] Solon Barocas, Moritz Hardt, and Arvind Narayanan. Introduction, 2021. Accessed: 2025-04-07.
- [6] Jan Batzner, Volker Stocker, Stefan Schmid, and Gjergji Kasneci. Germanpartiesqa: Benchmarking commercial large language models for political bias and sycophancy, 2024.
- [7] Django Beatty, Kritsada Masanthia, Teepakorn Kaphol, and Niphan Sethi. Revealing hidden bias in ai: Lessons from large language models, 2024.
- [8] Tommaso Del Becaro. Generative artificial intelligence and gender biases: Between new tools and human rights. *Student’s Social Science Journal*, 2018. Faculty of Law, University Goce Delchev Shtip.
- [9] Ruha Benjamin. *Race After Technology*. John Wiley and Sons, 2019.
- [10] Abeba Birhane and Fred Cummins. Algorithmic injustices: Towards a relational ethics. *arXiv preprint arXiv:1912.07376*, 2020.
- [11] Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (NeurIPS)*, pages 4349–4357, 2016.
- [12] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Proceedings of the 2018 Conference on Fairness, Accountability, and Transparency (FAT\*)*, volume 81 of *Proceedings of Machine Learning Research*, pages 1–15. PMLR, 2018.
- [13] Vanessa Cecil. *Older Women Navigating Age Stigma: Strategies and Outcomes*. Ph.d. dissertation, University of Exeter, 2024. ProQuest Document ID: 31876678.
- [14] Neil D. Christiansen and Caren B. Goldberg. Stereotypes at work: Occupational stereotypes predict race and gender distributions in the labor force. *Journal of Vocational Behavior*, 115:103318, 2019.
- [15] Juliet Corbin and Anselm Strauss. *Basics of Qualitative Research*. SAGE Publications, 2014.
- [16] Kate Crawford. *Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. Yale University Press, New Haven, 2021.
- [17] Amit Datta, Michael Carl Tschantz, and Anupam Datta. Automated experiments on ad

- privacy settings: A tale of opacity, choice, and discrimination. *Proceedings on Privacy Enhancing Technologies*, 2015(1):92,112, apr 2015.
- [18] Mara Del Baldo, Maria Gabriella Baldarelli, and Cristiana Ferrone. Gender diversity in accounting organizations: Empirical evidence and challenges in europe. *International Journal of Business and Management*, 14(11):1–15, 2019.
- [19] Artem Domnich and Gholamreza Anbarjafari. Responsible AI: gender bias assessment in emotion recognition. *CoRR*, abs/2103.11436, 2021.
- [20] Yogesh K. Dwivedi, Nir Kshetri, Laurie Hughes, Emma Louise Slade, Anand Jeyaraj, Arpan Kumar Kar, Abdullah M. Baabdullah, Alex Koohang, Vishnupriya Raghavan, Manju Ahuja, Hanaa Albanna, Mousa Ahmad Albashrawi, Adil S. Al-Busaidi, Janarthanan Balakrishnan, Yves Barlette, Sriparna Basu, Indranil Bose, Laurence Brooks, Dimitrios Buhalis, Lemuria Carter, Soumyadeb Chowdhury, Tom Crick, Scott W. Cunningham, Gareth H. Davies, Robert M. Davison, Rahul De, Denis Dennehy, Yanqing Duan, Rameshwar Dubey, Rohita Dwivedi, John S. Edwards, Carlos Flaviano, Robin Gauld, Varun Grover, Mei-Chih Hu, Marijn Janssen, Paul Jones, Iris Junglas, Sangeeta Khorana, Sascha Kraus, Kai R. Larsen, Paul Latreille, Sven Laumer, F. Tegwen Malik, Abbas Mardani, Marcello Mariani, Sunil Mithas, Emmanuel Mogaji, Jeretta Horn Nord, Siobhan O’Connor, Fevzi Okumus, Margherita Pagani, Neeraj Pandey, Savvas Papagiannidis, Ilias O. Pappas, Nishith Pathak, Jan Pries-Heje, Ramakrishnan Raman, Nripendra P. Rana, SvenVolker Rehm, Samuel RibeiroNavarrete, Alexander Richter, Frantz Rowe, Suprateek Sarker, Bernd Carsten Stahl, Manoj Kumar Tiwari, Wil van der Aalst, Viswanath Venkatesh, Giampaolo Viglia, Michael Wade, Paul Walton, Jochen Wirtz, and Ryan Wright. Opinion paper: "so what if chatgpt wrote it?" multidisciplinary perspectives on opportunities, challenges and implications of generative conversational ai for research, practice and policy. *International Journal of Information Management*, 71:102642, aug 2023.
- [21] Virginia Eubanks. *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor*. St. Martin’s Press, 2018.
- [22] Agneta Fischer and Marianne LaFrance. What drives the smile and the tear: Why women are more emotionally expressive than men. 7(1):22–29, dec 2014.
- [23] JihaD FraiJ and Varallyai Laszlo. A literature review: Artificial intelligence impact on the recruitment process. *International Journal of Engineering and Management Sciences*, 6(1):108–119, May 2021.
- [24] Ingrid Galster. *Le deuxième sexe de Simone de Beauvoir*. Presses Paris Sorbonne, 2004.
- [25] Estrella Gomez-Herrera and Sabine Koeszegi. A gender perspective on artificial intelligence and jobs: the vicious cycle of digital inequality. *Bruegel Working Paper*, (15), 2022. Bruegel Working Paper 15/2022.
- [26] UC Berkeley Haas School of Business. The playbook for addressing bias in artificial intelligence, 2021.
- [27] Hewlett Packard Enterprise. Cos’è il word embedding?, 2021. Accessed: 2025-04-18.
- [28] ICAEW Insights. Gender parity in the accounting profession: Still a long way to go, 2021. Accessed: 2025-04-01.
- [29] Elsevier Inc. Gender segregation an overview, n.d. Accessed via ScienceDirect Topics.
- [30] International Labour Office. International standard classification of occupations: Isco-08. volume i: Structure, group definitions and correspondence tables, 2012.
- [31] Matthew Kay, Cynthia Matuszek, and Sean A. Munson. Unequal representation and gender stereotypes in image search results for occupations. In *Proceedings of the ACM CHI Conference on Human Factors in Computing Systems*, Seoul, Korea, apr 2015.
- [32] Anna Sung Kelvin Leong. Gender stereotypes in artificial intelligence within the accounting profession using large language models. *Humanities and Social Sciences Communications*, 11(1), Sep 2024.

- [33] Yingqi Li and Ruihua Lu. When neutrality conceals bias: Perceived discrimination in algorithmic decisions. *European Journal of Social Psychology*, aug 2025.
- [34] Kirsten Martin. *Ethics of Data and Analytics*. CRC Press, 2022.
- [35] Joe Slater Michael Townsen Hicks, James Humphries. Chatgpt is bullshit. *Ethics and Information Technology*, 26(2), Jun 2024.
- [36] Kirsten Morehouse, Weiwei Pan, Juan Manuel Contreras, and Mahzarin R. Banaji. Bias transmission in large language models: Evidence from gender-occupation bias in GPT-4. In *ICML 2024 Next Generation of AI Safety Workshop*, 2024.
- [37] Silvia Moscatelli, Michela Menegatti, Naomi Ellemers, Marco Giovanni Mariani, and Monica Rubini. Men should be competent, women should have it all: Multiple criteria in the evaluation of female job candidates. *Sex Roles*, 83(5-6):269–288, jan 2020.
- [38] Palanichamy Naveen. The rise of ai in job applications: a generative adversarial tug-of-war. *AI and SOCIETY*, Aug 2024.
- [39] Safiya Umoja Noble. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press, New York, 2018.
- [40] Cathy O’Neil. *Weapons of Math Destruction*. Crown Publishing Group (NY), 2016.
- [41] Sara Panerati, Monica Rubini, Valeria A. Giannella, Michela Menegatti, and Silvia Moscatelli. A multidimensional implicit approach to gender stereotypes. *Frontiers in Psychology*, 14, nov 2023.
- [42] Marcelo OR Prates, Pedro H Avelar, and Luis C Lamb. Assessing gender bias in machine translation: A case study with google translate, 2019. Available from ResearchGate.
- [43] Manish Raghavan, Solon Barocas, Jon Kleinberg, and Karen Levy. Mitigating bias in algorithmic employment screening: Evaluating claims and practices. *SSRN Electronic Journal*, 2019.
- [44] Marco Rondina, Fabiana Vinci, Antonio Vetro, and Juan Carlos De Martin. Facial analysis systems and down syndrome. In *Machine Learning and Principles and Practice of Knowledge Discovery in Databases*, pages 145–160, feb 2025. First Online: 08 February 2025.
- [45] David Roselli, Kenneth R. Matthews, and Roger Varshney. Managing bias in ai. *IBM Journal of Research and Development*, 63(4/5):3:1–3:10, 2019.
- [46] Johnny Saldana. *The Coding Manual for Qualitative Researchers*. SAGE, 2021.
- [47] Helen Liu Sinead O’Connor. Gender bias—perpetuation and mitigation in ai technologies: challenges and opportunities. *AI and SOCIETY*, 39(4):2045–2057, Aug 2024.
- [48] Susan Sontag. The double standard of aging. In *The Other Within Us*, page 6. Routledge, London, 1st edition, 1997. Reprinted essay by Susan Sontag, originally published in 1972.
- [49] Yvonne Stedham, Jeanne H. Yamamura, and Hajime Satoh. Gender and leadership style: A comparison of u.s. and japanese accounting firms. *International Journal of Human Resource Management*, 17(5):660–676, 2006.
- [50] Luhang Sun, Mian Wei, Yibing Sun, Yoo Ji Suh, Liwei Shen, and Sijia Yang. Smiling women pitching down: auditing representational and presentational gender biases in image-generative ai. *Journal of Computer-Mediated Communication*, 29(1), Nov 2023.
- [51] Emeric Kubiak Tales Marra. Addressing diversity in hiring procedures: a generative adversarial network approach. *AI and Ethics*, May 2024.
- [52] Suze Wilson Toby Newstead, Bronwyn Eager. How ai can perpetuate – or help mitigate – gender bias in leadership. *Organizational Dynamics*, 52(4):100998, Oct 2023.
- [53] U.S. Department of Labor. O\*net online. <https://www.onetonline.org/>, 2025. Accessed: 2025-04-10.
- [54] Meghna Vidwans and Jeffrey R. Cohen. Gender disparities in accounting: Evidence from new zealand. *Pacific Accounting Review*, 32(4):501–520, 2020.
- [55] Judy Wajcman. *Feminism Confronts Technology*. Penn State University Press, University Park, PA, 1991.



- [56] Judy Wajcman. *TechnoFeminism*. John Wiley and Sons, 2013.
- [57] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5310–5319, 2019.
- [58] Wikipedia contributors. Generative artificial intelligence — Wikipedia, the free encyclopedia, 2025. [Online; accessed 3-April-2025].
- [59] Mi Zhou, Vibhanshu Abhishek, Timothy Derdenger, Jaymo Kim, and Kannan Srinivasan. Bias in generative ai, 2024.