

“To Benefit all Humanity”: Towards Fair Algorithmic Systems through Substantive Equality and Theories of Justice

MASTER THESIS

to obtain the Erasmus Mundus Joint Master Degree
in Digital Communication Leadership (DCLead)

of

Faculty of Cultural and Social Sciences
Paris Lodron University of Salzburg

Technical Faculty of IT and Design
Aalborg University in Copenhagen

Submitted by
John Gustavo Choque Condori

20216817

hello@johnchque.com

Rebæk Søpark 5, 2650 Hvidovre, Denmark

Jannick Kirk Sørensen (Aalborg University Copenhagen)

Josef Trappel (Paris-Lodron Universität Salzburg)

Department of Communication Studies

Copenhagen, 31 July 2023

Table of Contents

1. Introduction.....	7
1.1. Background.....	7
1.2. Research Questions.....	8
1.4. Relevance.....	9
1.3. Methodology.....	10
1.5. Scope and Limitations.....	12
1.6. Structure.....	12
2. Literature Review.....	14
2.1. Literature Review Process.....	15
2.2. Fairness in Psychometrics: a Paradigm Problem.....	16
2.3. Fairness in Computer Systems: a Positivist Approach.....	19
2.3.1. Fairness in complex systems.....	19
2.3.2. Network fairness.....	20
2.4. Fairness in Machine Learning: a Substantial Paradigm Problem.....	22
2.4.1. Harms and consequences of machine learning models.....	24
2.4.2. Incompatible fairness notions.....	26
2.4.3. The impossibility of algorithmic fairness.....	28
2.4.4. Escaping the impossibility of fairness.....	29
3. The Wickedness of Algorithmic Fairness.....	31
3.1. Wicked Problems Framework.....	31
3.2. Algorithmic Fairness is a Wicked Problem.....	34
4. Theoretical Framework: Shifting Focus.....	42
4.1. Theories of Justice.....	43
4.1.1. Justice as fairness.....	44
4.1.2. The capabilities approach.....	46
4.2. Substantive Algorithmic Fairness.....	49
5. Analytical Framework: A Proposal for Algorithmic Justice.....	52

5.1. Step 1: Diagnose inequalities.....	54
5.2. Step 2: Identify potential reforms.....	54
5.3. Step 3: Consider roles of algorithms.....	56
6. Experimental Setup: Beyond Fairness Metrics.....	58
6.1. Data Collection.....	59
6.1.1. Dataset.....	60
6.1.2. Models.....	61
6.2. Data Analysis.....	62
6.3. Ethical Considerations.....	64
7. Results: Identifying social inequalities.....	65
7.1. Accuracy and Bias: Quantitative results.....	65
7.2. Framework Analysis: Qualitative Results.....	69
8. Discussion: Towards Algorithmic Justice.....	73
9. Conclusion: “To Benefit all Humanity”	77
9.1. Future research.....	80
10. References.....	83
11. Appendix.....	96
11.1. Appendix A: Transformer model architecture.....	96

List of Figures

Figure 1. Substantive Algorithmic Fairness framework (Green, 2022).....	51
Figure 2. A novel framework for Algorithmic Justice. Adapted from Green (2022).....	53
Figure 3. Model accuracy using the original dataset.....	67
Figure 4. Model accuracy using the alternative dataset.....	67
Figure 5. Model bias score using the original dataset.....	68
Figure 6. Model bias score using the alternative dataset.....	68

List of Tables

Table 1. List of definitions.....	43
Table 2. List of transformer models used.....	59
Table 3. Prompts with wrong outputs in more than one question variation.....	69

*To my family and friends,
For being the light piercing through the shadows.
To the disadvantaged,
For you, the sun will rise in tomorrow's skies.*

Thank you.

Executive Summary

This thesis explores the challenges of algorithmic fairness and proposes a novel framework for developing fair algorithmic systems based on philosophical theories of justice. A historical view of the problem is presented, considering the foundation of fairness research in computer systems, related problems in psychometrics, and an elaborated view of the current issues in algorithmic fairness. The research establishes algorithmic fairness as a wicked problem due to the lack of a widely accepted definition, resulting in conflicting mathematical fairness metrics. To address this, the study reframes the discussion around algorithmic justice, which expands the scope of analysis beyond specific decision points and considers the voices of the ones impacted by algorithmic systems. Moreover, the Substantive algorithmic fairness framework is introduced as a means to promote justice in practice and address social hierarchies in algorithmic decision-making. Building upon this framework, the thesis presents a novel analytical approach that connects relational and structural considerations with representational and allocative harms. By using theories of justice, potential reforms can be identified to create fair algorithmic systems. Through experiments and qualitative analysis, the proposed approach is applied to transformer models, revealing biases related to multiple categories and inconsistencies in the model outputs. The study concludes that integrating philosophical theories of justice can lead to fair algorithmic systems that align with societal principles and benefit humanity as a whole, offering new avenues for addressing the wickedness of algorithmic fairness.

1. Introduction

1.1. Background

“Not everything that counts can be counted, and not everything that can be counted counts” - William Bruce Cameron

Being educated as a computer scientist, most of my understanding of the world was based solely on mathematical representations and data flows. I spent the early days of my professional journey developing open source tools serving thousands of users around the world. The more I learned about the industry, the more I understood about the positive impact that technology can have in society. However, throughout this journey, I found myself experiencing digital disparities along multiple fronts. From the lack of support for my native language in prevalent tools, the inaccurate facial recognition technologies for darker skin tones (BBC News, 2015), to the public experimentation of Facebook algorithms in my region (Cellan-Jones, 2017).

Throughout my Master’s education, I eventually realised that the reasons for these digital disparities were rooted in a variety of social and technical factors. This led me to develop a strong interest in understanding how the latest developments in machine learning models could address these complexities. However, while keeping up with the latest developments, I soon realised that these models have been shown to reproduce social biases (Dehouche, 2021; Kirk et al., 2021; Ray, 2023). Because machine learning models are trained with labelled data from the real world, they can reflect biases that arise from underlying patterns of discrimination (Binns, 2018; Salimi et al., 2020). These models will also encode the biases of their developers and make assumptions in order to make real-world decisions (Baker & Hawn, 2022; Friedler et al., 2021; M. S. A. Lee & Floridi, 2021). This exacerbates the problem, as a biased model can perpetuate discrimination at scale (M. S. A. Lee & Floridi, 2021), they can be discriminatory against minorities in multiple ways (Fletcher et al., 2021; Johnson et al., 2022; Paulus & Kent, 2020), and can be used for misinformation, spam, phishing, fraudulent academic essays, and more (Dehouche, 2021).

Even though several fairness notions have been proposed in the literature, multiple authors claim that these proposed measures are mathematically incompatible with each

other (Ashokan & Haas, 2021; Castelnovo et al., 2022; Green, 2022; Verma & Rubin, 2018). This incompatibility between notions raises what is called the “impossibility of algorithmic fairness” (Green, 2022), meaning that it is impossible for an algorithm to satisfy all fairness notions at once. Such a problem was previously observed in the field of psychometrics, when the authors were unable to set a widely accepted definition of the concept, leading to the development of multiple statistical models that were incompatible with each other (Cole & Zieky, 2001). This can be rooted in the positivist nature in which fairness has been studied in computer systems, such as the bandwidth allocation of networks and by preventing infinite computations in complex systems.

This raises multiple questions about the operationalization of what it means to be fair or non-discriminatory (Binns, 2018). The existing fairness definitions seem to overlook the ethical and philosophical definitions of the concept (Card & Smith, 2020; Green, 2022), while the industry of AI development lacks ethical training and applications (Munn, 2022). As a consequence, in order to escape this mathematical incompatibility, multiple authors have defined a set of principles to achieve algorithmic fairness by defining what is expected from this mathematical notions (Giovanola & Tiribelli, 2022; Paulus & Kent, 2020; Zehlike et al., 2020). However, some authors such as Munn (2022) argue that in order to develop fair algorithmic systems, alternative approaches should go beyond the definition of ethical principles and consider more broadly the systems of oppression. In this way, AI justice reframes the discussion around AI ethics to consider that the moral properties of algorithms are not internal to the models and instead are a product of the social systems where they operate (Gabriel, 2022).

By undertaking this research, this thesis aims to propose a different perspective to the discussion of algorithmic fairness by considering theories from philosophy, and to understand the current challenges that would prevent its achievement. The complexity of the topic may lead to more questions than answers. However, with this in mind, I hope that this work will be memorable for my future endeavours, while also being a modest contribution to the development of fair algorithmic systems.

1.2. Research Questions

- *What factors influence the effective operationalization of fairness for the implementation in algorithmic systems?*

- *How can theories of justice be used to develop a framework for implementing fair algorithmic systems that effectively addresses and mitigates the challenges posed by the wickedness of algorithmic fairness?*

1.4. Relevance

In 2017, Vaswani et al. (2017) published the paper *Attention is all you need*. In the paper, the authors introduce the transformer model architecture, which uses the attention mechanism as the only mechanism with which to derive dependencies between inputs and outputs in language models (Amatriain, 2023). Since then, transformer models have evolved at an incredible rate. This revolution achieved an important milestone when OpenAI announced the release of ChatGPT (OpenAI, 2022). ChatGPT was presented as a conversational AI that can chat, answer questions, and challenge incorrect assumptions (OpenAI, 2023b). After 5 days of its release, ChatGPT surpassed Instagram to become the quickest application to reach 1 million users (UBS Editorial Team, 2023). This trend was followed by other companies when they introduced their own AI chatbots to challenge ChatGPT. Google announced the release of Bard (Pichai, 2023). In the same month, Alibaba, JD.com, and Tencent from China announced that they were working on their own AI products to challenge ChatGPT (Sorkin et al., 2023).

While more companies embraced the new AI race and announced the development of their own solutions, ChatGPT was found to present multiple ethical issues. In December 2022, a Twitter user reported that while testing the model with code generations, ChatGPT suggested that if an individual was from North Korea, Syria, or Iran, they should be tortured (Biddle, 2022). During the same months, ChatGPT had already been found creating text that implied that only white and asian men would be good scientists and that a woman in a lab coat would only be there to clean the floor (Alba, 2022). Some other authors reported that, while the model would not generate a racist story from a prompt that asks for it, asking the model to behave like a racist writer would generate such a story (Vock, 2022). Soon after, similar claims were made for Bard as it was reported that it presented liberal inclinations (Durden, 2023; Reinl, 2023).

As the transformer revolution has allowed the development of ChatGPT and Bard to become widely popular, multiple ethical discussions have arisen. Language models may reflect and amplify stereotypes (Dehouche, 2021; Kirk et al., 2021; Ray, 2023) or reflect

the biases from society and their developers (Baker & Hawn, 2022; M. S. A. Lee & Floridi, 2021). This has caused an increasing concern over the accountability and regulation of these models. However, as previously mentioned, multiple fairness notions have been defined and often been mathematically incompatible with each other. This raises multiple questions about the operationalization of what it means to be fair or non-discriminatory (Binns, 2018). This work focuses on this research gap and seeks to make a contribution in exploring alternative ways to implement fair algorithmic systems.

1.3. Methodology

In order to answer the research questions of this thesis, two main sections are defined. First, the concept of fairness is explored from different perspectives, following an inductive approach. This led to the argument that algorithmic fairness is a wicked problem. Second, the study explores theories of justice and other frameworks to propose a novel approach in order to escape the wickedness of algorithmic fairness.

In Chapter 2, an inductive literature review is presented. It aims to gain a comprehensive understanding of the various perspectives of fairness in different domains. It describes the multiple struggles to identify a widely accepted definition of fairness in psychometrics, due to the contextual nature of the term. At the same time, it describes the foundation of the study of fairness in computer systems, by exploring fairness in networks and complex systems. With this understanding, a more critical review of the literature of algorithmic fairness is presented, showing the multiple problems caused by the lack of a widely accepted definition of the term. This draws the theoretical foundation to argue that fairness is a contested topic that mirrors the on-going debate between positivist and constructivist approaches for research. As computer systems set a positivist foundation for the study of fairness in computer systems, algorithmic fairness is often described in mathematical models and fails to account for the systematic and historical issues of inequalities.

Building upon the insights from the literature review, the wicked problem framework is used as a critical lens to examine the concept of algorithmic fairness. Throughout Chapter 3, multiple arguments from the literature are presented showcasing the diverse set of fundamental issues that the study of fairness in algorithmic systems present. The

chapter concludes with the critical issues that a novel approach should consider in order to escape the wickedness of the problem.

To address the wickedness of algorithmic fairness, the rest of this work follows an exploratory approach. It begins by defining a clear differentiation between fairness and justice, as the latter is the preferred term in the literature of philosophy. A strong theoretical foundation enabled the author to delve into other research made around AI justice. In consequence, a novel framework for algorithmic fairness was defined through an inductive approach that involved connecting the Substantive algorithmic fairness framework by Green (2022) with theories of justice and the representational harms proposed by Crawford (2017). By combining these diverse perspectives, the study sought to develop a comprehensive and nuanced understanding of algorithmic fairness, taking into account both ethical considerations and the potential impact of representational biases.

The author conducted an experiment to assess the representational biases present in transformer models. The dataset used for the experiment was selected based on its applicability to transformer models and its theoretical foundation on representational harms. The experiment was conducted with two different models, GPT-3.5 and Alpaca LoRA, and their accuracy and bias scores were calculated to quantify the biasness in their outputs. The study continues with a qualitative analysis of the outputs. By applying the capabilities proposed by Nussbaum (2011), the research aimed to gain insights into the alignment of the models with these capabilities. This qualitative examination provides more context to understand how the models may impact the lives of individuals and whether they currently uphold principles of justice.

By employing an interdisciplinary methodology, combining theories of philosophy with algorithmic systems, the research aimed to demonstrate the potential of applying philosophical theories of justice in the analysis and implementation of fair algorithmic systems. The combination of the novel framework, the experimentation with transformer models, and the qualitative analysis allowed for a multi-faceted exploration of fairness, contributing to the ongoing discourse surrounding algorithmic fairness.

1.5. Scope and Limitations

The aim of this thesis is not to have an exhaustive list of fairness notions or to map the relationship between them. Such work has been made in previous studies (Ashokan & Haas, 2021; Castelnovo et al., 2022; Green, 2022; Mehrabi et al., 2021; Verma & Rubin, 2018). Instead, this thesis aims at highlighting the challenges of implementing fairness in algorithmic systems and to explore novel approaches for such a task. Further comments about possible future work are described in the conclusion.

While the concept of fairness has been discussed in a variety of fields (Duran-Rodas et al., 2020; Garay et al., 2006; Hamermesh & Schmidt, 2003; Howard et al., 2016; Krysiak & Krysiak, 2006; L. Li et al., 2014; Moret et al., 2020; Persico, 2002) this work does not aim to comprehensively cover all domains in which fairness is discussed. The scope of this thesis is limited by time and resource constraints thus, this work considers the discussion of fairness in machine learning, which represents the most extensively studied algorithmic system. Additionally, this study explores fairness in complex systems and networks, as they form the foundation for understanding fairness in computer systems. Furthermore, the discussion of fairness in psychometrics is included due to its numerous similarities with the ongoing discussions in machine learning models.

The discussion of algorithmic fairness opens up multiple questions about the applicability of fairness in real world settings. This master thesis does not aim to find answers to these questions or to propose a novel fairness metric general enough to expand its applicability to multiple contexts. Instead this master thesis is a small step into analysing a different direction by discussing the impossibility of algorithmic fairness and placing philosophical theories as an important part of the discussion. This approach does not disregard the previous research efforts related to algorithmic fairness and instead extends them to consider the systematic social problems that lead to unfairness in algorithmic systems.

1.6. Structure

The following chapters are structured as follows. Chapter 2 presents a comprehensive literature review in the domains of psychometrics, computer systems, and machine

learning. This chapter highlights the main findings, gaps, and challenges related to algorithmic fairness and it serves as the foundation for the subsequent chapters. Chapter 3 uses the knowledge gathered in Chapter 2 to argue that algorithmic fairness is a wicked problem, resulting in the need for a shift in focus for its implementation. It provides a detailed discussion of the wickedness of algorithmic fairness and highlights the considerations for newer approaches beyond fairness metrics.

Chapter 4 defines a clear distinction between the terms fairness and justice and provides a thorough description of two of the most prominent modern theories of justice, serving as a theoretical foundation for a framework of algorithmic justice. The chapter also introduces the Substantive algorithmic fairness framework as a new way of understanding fairness and promoting justice in practice. Chapter 5 introduces a novel approach to applying the Substantive algorithmic fairness framework and the previously mentioned theories of justice in transformer models. This chapter explains the iterative process of promoting justice in practice by identifying relational and structural harms, defining potential social reforms, and the role of algorithms to implement these reforms.

Chapter 6 presents an experiment that applies the framework of Chapter 5 to transformer models. The experiment focuses on the relational perspective of the proposed framework by using a dataset constructed to identify biases associated with representational harms. The chapter describes the experimental design, data collection, and data analysis methods. It provides details of the dataset and the transformer models selected, along with the criteria for their selection. Chapter 7 provides the calculation of accuracy and bias scores for the selected transformer models and it describes the findings of the qualitative analysis of the outputs by using the capabilities approach. Chapter 8 follows with a comprehensive discussion of the findings and the considerations from all previous chapters to answer the research questions. It expands on the nature of the discussion, highlights the key insights and implications, and addresses the research questions directly to provide a conclusion to the analysis. Finally, chapter 9 and 10 include the conclusion for the thesis, a summary of the main contributions, future research directions, and the references used throughout the thesis.

2. Literature Review

The discussion of developing fair algorithmic systems has been approached from different perspectives, reflecting the ongoing debate between positivist and constructivist approaches to research. Some authors have proposed mathematical notions of fairness to quantify and measure the fairness of algorithmic systems, while others have focused on discussing the implications for society and the potential harms that unfair systems can cause. This duality of perspectives highlights the need to examine the challenges and complexities of achieving fairness in practice in algorithmic systems.

To understand these challenges, this literature review will provide some context by describing the challenges faced in developing fair tests in psychometrics. Initially conceived as a purely mathematical problem, it was later recognized that decisions made at specific points may fail to reflect existing patterns of discrimination, and that maximising statistical models alone would not take these issues into account. Building upon this foundation, this chapter will then move on to a discussion of fairness in computer systems, particularly in the context of complex systems and networks. This section will explore how fairness has been predominantly studied following a positivist approach in the foundations of modern computing.

The chapter concludes with a thorough review of the discussion of fairness in machine learning models. This section will reveal that the literature on fairness in machine learning follows a similar pattern to the studies in psychometrics. Its foundation is rooted in incompatible mathematical models that do not account for the potential patterns of discrimination in society. This issue becomes even more complex with the implementation of machine learning models, which have the potential for large-scale implementations that may reproduce the patterns in their data and have a negative impact on disadvantaged groups. This section will describe the potential harms that unfair systems may cause, conflicting fairness notions proposed in the literature, what is considered the impossibility of algorithmic fairness, and what other authors have tried to escape this impossibility.

2.1. Literature Review Process

In order to provide a solid theoretical basis for the discussion of algorithmic fairness, it was important to start by reviewing how fairness has been, defined, operationalised, and discussed in previous studies. This meant to use an inductive approach with the literature review. To do so, a database search was performed using the Scopus database. The initial query used the term *fairness* to search occurrences in the article title, abstract, and keyword. The result from such a query returned more than 60 thousand papers in nearly 30 different fields. The vast amount of results required to try different approaches to get a smaller number of results that were reviewable. Using fairness as a keyword also resulted in more than 9 thousand papers found. Limiting these results by relevance and dates was avoided in order to have historical development of the term and not relying on the metrics defined by the database used. After multiple iterations, the search query used for the present literature review was defined as:

```
TITLE-ABS-KEY ( "fairness" W/2 "definition" ) AND ( LIMIT-TO ( PUBSTAGE ,  
"final" ) ) AND ( LIMIT-TO ( DOCTYPE , "ar" ) ) AND ( LIMIT-TO ( LANGUAGE ,  
"English" ) )
```

With this query, the term fairness and definition should be separated by 0 to 2 words. This decision was taken as multiple papers found in the original queries showed definitions using other words in between, such as *definition of bandwidth fairness*. At the same time, the query specifies only research articles written in English. Any other combination of synonyms to *definition*, such as *meaning* or *notions* were not included as they caused a large increase in the results. Moreover, in order to focus on understanding the definition of fairness, no similar terms were used, such as *justice* or *bias*. Still, some papers that discussed both were returned in the final result from the query above.

This query returned 142 papers from which an initial screening of their abstracts was performed. Two exclusion criteria were considered: 1) papers that focused on mathematically testing previously defined definitions, and 2) papers that did not discuss or define fairness. During this process, duplicates were also removed leading to end up with 93 relevant papers. The literature review is presented below.

2.2. Fairness in Psychometrics: a Paradigm Problem

Within the literature that discusses fairness, psychometrics is one of the earliest fields that covered the topic extensively. The main purpose of psychometrics is to explain the meaning of responses given by subjects and to suggest techniques for measuring mental processes (Pasquali, 2009). Different studies argue that fairness does not have an agreed-upon definition that is widely accepted causing a variety of statistical models to be proposed in order to clarify its broad nature (Baharloo, 2013). Such an issue was already acknowledged by Darlington (1971) who argues that the term *fair* carries conflicting connotations so that no single measurement test can meet all specifications required to make a *culturally fair* test. In consequence, fairness may be defined depending on what weights (importance) are allocated to specific criteria by decision-makers (Steffy & Ledvinka, 1989). These decision-makers may consciously or unconsciously use a definition of selection bias for a selection process (Ledvinka, 1979). As well as acknowledging that the same test outcome may be seen positively or negatively depending on the value of different groups (Cole & Zieky, 2001).

The theory of fairness was shown as a complex issue that had real-life implications in the fair treatment of minority groups and their selection (Flaughner, 1974). Baron (2017) describes that most people define fairness in one of three ways: 1) Distributive justice, where the contributions that individuals make are reflected in the rewards they receive; 2) All people should receive the same regardless of how much they contribute; and 3) Personal needs, meaning that fairness exist when rewards are divided as needed by individuals. However, fairness may not be considered a concern to the people that believe that the differences between groups are just a reflection of reality (Baharloo, 2013).

According to Cole & Zieky (2001), during World War I, educators started taking advantage of large-scale testing and group differences were thought to be a reflection of reality. However, by 1960, more professionals paid increasing attention to the difference in scores between groups and focused on test and item fairness. This increasing interest was also influenced by the publishing of the paper *How Much Can We Boost IQ And Scholastic Achievement?* (Jensen, 1969) arguing to consider a genetic component as a cause for differences in performance between white and african-american test takers.

From this discussion a line of research emerged to discuss the development of models of fair test use in selection processes. Despite illuminating the complexities of defining fairness, it failed to yield consensus in the meaning of fairness and widely implemented procedures (Cole & Zieky, 2001). As a consequence, multiple models were proposed, each one challenging the assumptions of the previous one and considering additional levels of assumptions. The Cleary model (Cleary, 1968) defines that a test is biased if in the prediction of an outcome, there are consistent errors between the actual outcome and the predicted outcome of the test for members of a subgroup. However, this model was highly criticised as it was shown that it maximises utility and minimises minority employment (Ledvinka, 1979; Steffy & Ledvinka, 1989).

Alternative models were later proposed that challenged the assumptions of the Cleary model. These include the Thorndike model (Cole & Zieky, 2001; Flaughner, 1974; Steffy & Ledvinka, 1989), the Cole model (Cole, 1973), and the Darlington model (Flaughner, 1974). However, all these models cannot be true at the same time, as these models are not mathematically compatible under all selection conditions (Cole & Zieky, 2001; Ledvinka, 1979). The development of these additional models for fairness were led by the focus on the consequences of test use in selection and how the difference between test scores and criterion scores displayed a potential biased criteria (Cole & Zieky, 2001). Moreover, these debates intensified when tests could be equally valid for different groups but unfair to one depending on the fairness of definition considered (Cole & Zieky, 2001).

By the 2000s, research on fairness had reached a more general consensus that the concept depends on different social circumstances. For instance, in the Standards for Educational and Psychological Testing from the United States, fairness was indicated to be subject to different definitions and interpretations in the different social circumstances (Cole & Zieky, 2001). Moreover Cole & Zieky (2001) argue that no statistic could indicate unambiguously if a test item is fair. This caused some education teachers to be reluctant to teach students with disabilities as their inclusion may be considered unfair to typically achieving students or that these inclusion may cause more demands on education teachers (Berry, 2008). These issues were related to whether fairness means differential treatment for disadvantaged groups or equal treatment to everybody (Zollers et al., 2000).

Students may refer to their personal expectations when discussing fairness instead of considering differential prediction, adverse impact, or personal bias (Crocker, 2003). The application of justice principles in grading, including equity, consistency, voice, and justification, can impact students' perception of grading fairness, and their perception of the teacher's adherence to these principles determines their perception of fairness (Rasooli et al., 2019). On the other hand, teachers may define their attitudes towards fairness as a manifestation of their personal philosophy of justice (Berry, 2008). In consequence, Baharloo (2013) suggests that a different approach to traditional assessment, dynamic assessment, would integrate instruction and assessment in all learning stages so they cannot be distinguishable along the process.

For Deutsch (1985), some notions of fairness derive from justice, meaning that they are concerned with individual well-being and societal functioning. This way, the values of justice are those that foster effective cooperation that promote well-being. However, these goods and conditions may be limited, so rules must be made to facilitate this distribution. Distributive justice creates rules to guide such distribution either by considering the characteristics and needs of recipients or by giving everyone equal benefits (Berry, 2008). These ideas are deeply related to two principles of Justice as fairness, a theory of justice presented in Chapter 4, where Rawls explains that unequal rules are justified if they favour the individuals in disadvantage. Similarly, Berry (2008) describes needs-based principles of distributive justice when the well-being of individuals is the main concern and giving priority to the ones less well off. This means that all students receive the support they need to thrive.

Multiple different definitions of fairness have been defined over the years, some of them focusing more on the disparate treatment of individuals and how morally wrong that is in itself, while others compare different groups and apply such knowledge to testing practices across groups (Baharloo, 2013; Berry, 2008). However, by the 2000s, Cole & Zieky (2001) argued that research hasn't supplied a generally accepted definition of fairness and there are no analyses that unequivocally indicate the existence of unfairness. After all, test fairness is a multi-faceted issue which goes beyond the content of the test and includes other aspects of testing as well (Baharloo, 2013). A dynamic assessment that adapts to the needs of every student requires more attention on the part of educational institutions (Baharloo, 2013). However, historically, the assessment

literature has theorised in large-scale traditions instead of adopting a classroom assessment perspective (Rasooli et al., 2019).

2.3. Fairness in Computer Systems: a Positivist Approach

While the discussion of fairness in psychometrics struggled to operationalise the concept while taking into account the perspectives of the affected groups, the early studies of fairness in computer systems were mostly concerned with the optimization of mathematical models to improve resource allocation. This section provides the necessary context for understanding the current discourse on fairness in algorithmic systems. It highlights the historical emphasis on mathematical optimisation and positivism in complex systems and networks, setting the stage for the modern discussion of fairness in computer systems.

2.3.1. Fairness in complex systems

Early literature discusses fairness in the context of complex systems. When non-deterministic models are used to represent complex systems, fairness becomes an important subject of study as the choice to take a course of action is not defined (Apt et al., 1988; Queille & Sifakis, 1983). Non-deterministic models are mathematical models that produce different outcomes at different times with the same input data (*Non-Deterministic*, 2023). These non-deterministic systems may include airline reservation systems, operating systems, concurrent algorithms, etc (Kwiatkowska, 1989). For instance, if three customers are on a waiting list to book a flight, if the first one tries to do a booking after a previous booking, the system might pick the second and third customers' booking processes after each other. If the second and third keep requesting bookings, the first might never be served (Kwiatkowska, 1989).

Multiple authors begin their discussion with a general definition of fairness, asserting that fairness guarantees that if a repeated nondeterministic choice is possible sufficiently often, then it proceeds sufficiently often (Apt et al., 1988; Queille & Sifakis, 1983; Wabenhurst, 2003). As the example above illustrates, the sequence of occurrences may lead to unfair situations where an event that can occur infinitely often does not because conflicts are not resolved in an equitable way (Queille & Sifakis, 1983). In other words, fairness notions are motivated by disallowing infinite computations so that

system components are not prevented from proceeding (Kwiatkowska, 1989). In this sense, while all finite computations are fair, it may be necessary to distinguish between fair and unfair infinite computations (Kwiatkowska, 1989).

Multiple observations are added to this definition. For instance, Apt et al. (1988) argue that depending on the definition of what “choice”, “possible”, and “sufficiently often” mean, multiple fairness notions may arise. This highlights a problem of interpretation and operationalization that was seen in the previous section related to psychometrics. For Queille & Sifakis (1983) it is important to consider that such a definition of fairness would only be fair to some set of states and not for others sets. Moreover, this intuitive definition considers every component in isolation and its properties are affected by the system that is taken into consideration, such as processes, events, transitions, or states (Kwiatkowska, 1989). In this context, fairness is then a mathematical abstraction that in multiprogramming environments, abstracts the details of fair schedulers, and in distributed environments, abstracts the speed of processors (Alur & Henzinger, 1998).

As a consequence, this characterization of fairness induces different notions corresponding to different events or possibilities that can occur (Queille & Sifakis, 1983). Apt et al. (1988) argue that the selection of a definition of fairness for a model relies almost exclusively on subjective criteria. Moreover, Kwiatkowska (1989) argued that the lack of consensus on the definition of fairness may come from the variety of semantic models and how fairness notions rely on the intrinsic characteristics of these models.

2.3.2. Network fairness

A different thread of research studied fairness in the context of bandwidth allocation. Zukerman & Chan (1993) explain that as ATM networks become more established and efficient, the resources are divided based on competition when the total demand is greater than the available resources. Thus, defining how the bandwidth will be shared in order not to discriminate against heavy users will require the definition of a fairness notion determining the proportions. In other words, fairness can be understood as the sufficient balance in the allocation of bandwidth for different entities (Bisio & Marchese, 2014). For example, if two satellites are unequally affected by the environment but they get the same bandwidth allocation, it would cause a performance unbalance between

the entities thus, treating them equally does not mean they are treated fairly (Bisio & Marchese, 2014). As a consequence, fairness is a key issue when discussing bandwidth allocation (Fei & Yang, 2005).

One of the best-known definitions of fairness for bandwidth allocation is Max-min fairness. It can be defined as “maximise the allocation of each user i subject to the constraint that an incremental increase in i ’s allocation does not cause a decrease of some other user’s allocation that is already as small as i ’s or smaller” (Chen et al., 1993). In general, comparing the relative utilities that every receiver in the network obtains can be used to measure the fairness of an allocation (Rubenstein et al., 2002). According to Nicosia et al. (2017), max-min fairness is based on the principles of justice of John Rawls in a way that even the least advantaged agent gains as much as possible. This insight reflects a similar analysis to that in psychometrics, where rules of distributive justice may be used to facilitate a class assessment approach. However, it was demonstrated that having max-min local fairness in a network does not translate to have max-min fairness in the network as a whole when insubordinate users are present (Chan & Zukerman, 2002). This shows that fairness is influenced by the context and level of analysis of a system.

Some solutions to this allocation problem include other factors. For example, Zukerman & Chan (1993) explain how in ATM networks, fairness is achieved when one user’s throughput, the amount of data transmitted in a given time, is not less than others in the same bottleneck so the capacity is shared equally between users. Alternatively, Chen et al. (1993) view the network as a collection of communication resources and suggest a custom algorithm for local fairness in order to ensure the use of network resources between nodes that compete for a subset of links. Other authors suggest fairness definitions that expand on max-min fairness. For instance, Active fairness refers to the use of dynamic weights adjusted to achieve the desired bandwidth allocation considering the network conditions to set them (Chew & Gupta, 2000). Moreover, more authors focus on the fairness of individuals, such as intra-session fairness (Fei & Yang, 2005), multipoint-to-multipoint fairness (J.-S. Li et al., 2007). Alternatively, others expand these definitions to newer technologies such as Max-min fairness for wireless and wireline networks (P. Wang et al., 2008), Multicast-favourable max-min fairness (Österberg & Zhang, 2011).

Since the early studies of bandwidth allocation, there has not been consensus on the definition of fairness in multipoint connections or between multipoint and point-to-point connections (Nguyen & Katzela, 2000). For Bisio & Marchese (2014), many research papers claim that their own fairness definition is fair even when it is defined in constraint to the definition of other papers. In the context of exchange protocols and satellite environment, other authors argue for the need of a unified view of fairness in order to avoid unambiguous interpretations (Bisio & Marchese, 2014; Markowitch et al., 2003). However, Fei & Yang (2005) investigate future heterogeneous fairness definitions as different receivers may want to use different functions. Moreover, J.-S. Li et al. (2007) argues that every autonomous system should have its own policy for fairness. In modern times, the discussion of fairness turned into the evaluation of the trade-off between fulfilling a fairness criterion and guaranteeing an adequate level of efficiency (Nicosia et al., 2017).

2.4. Fairness in Machine Learning: a Substantial Paradigm Problem

The previous sections provided a historical view of the studies of fairness in psychometrics and computer systems. This section will use the presented foundation to describe the predominant mathematical approach towards algorithmic fairness and how modern studies request for a broader level of analysis to account for the systematic discrimination patterns. This section highlights the challenges posed by the duality of approaches to study fairness in algorithmic systems, drawing parallels with the discussion in psychometrics. It examines the ongoing discourse that recognises the inherent impossibility of achieving algorithmic fairness due to the mathematical incompatibility of metrics. Finally, it examines alternative approaches proposed by scholars to overcome this impossibility and mitigate unfair outcomes.

Most of the research of algorithmic fairness has focused on developing ways to measure it and how to mitigate biases (Castelnovo et al., 2022). However, Card & Smith (2020) argue that a vast amount of this research has been made without an elaborated understanding of ethical foundations. In consequence, the literature also reflects that machine learning models will reflect in their outputs the historical biases and social inequalities with which they were trained (Salimi et al., 2020). Moreover, the models

will encode the biases of the developers that created them and end up being discriminatory against specific groups (Baker & Hawn, 2022; M. S. A. Lee & Floridi, 2021).

Machine learning models are currently implemented in a variety of fields. These include determining criminal sentences, hiring, loan applications, dropout prediction, essay scoring, graduate admission, among multiple other fields (Baker & Hawn, 2022; Friedler et al., 2021). However, in fields such as credit scoring, research shows a lack of real-world data available as institutes are not willing to share the data they gather (Szepannek & Lübke, 2021). Moreover, due to the issues of such data availability, computational limitations, and interpretability, research focuses on simpler models (Scutari et al., 2022). For instance, in credit scoring, logistic regression is still the most accepted method for classification even when modern algorithms have been demonstrated to have more potential benefits (Szepannek & Lübke, 2021).

The goal of implementing fair models is to prevent discrimination of individuals or groups following their characteristics (Aler Tubella et al., 2022). Such an idea can be considered a natural consequence of the promise that artificial intelligence is able to worsen human biases and attenuate disparities (Paulus & Kent, 2020). In that sense, fairness is generally understood as the concept that all groups and users are treated equally (Ashokan & Haas, 2021). However, the literature defines multiple mathematical fairness notions that were proved to be incompatible with each other (Ashokan & Haas, 2021; Castelnovo et al., 2022; Green, 2022; Verma & Rubin, 2018). A more thorough explanation of this incompatibility is presented in the following sections. In consequence, scholars have called for rejecting the idea of fairness and replacing it with justice, equity, or reparation (Green, 2022). Since these mathematical formulations overlook the contextual and philosophical meaning of fairness (Green, 2022).

Bias can be defined in multiple ways. For N. T. Lee et al. (2019), biases are the outcomes that are less favourable to individuals within groups where groups have no relevant differences that would justify such harms. Alternatively, it can also be used to describe when an algorithm has an unwanted dependence on an attribute in the data that belongs to a demographic group (Fletcher et al., 2021). Fletcher et al. (2021) argue that even if bias is related to fairness, it is independent from ethics and is only a

mathematical consequence of an algorithm and the input data. However, while developing these machine learning models, would require specific modelling choices that would create a natural trade-off of engineering and ethical considerations (Aler Tubella et al., 2022). As a consequence, multiple debiasing techniques have been developed and researched to equalise the outcomes to match the expectations of the multiple existing fairness notions (Aler Tubella et al., 2022).

If a model is found to contain bias, it can be judged against ethical and legal principles that would allow the tuning of the model to satisfy fairness criteria (Fletcher et al., 2021). This tuning may happen in different instances of the implementation of models, such as: data collection, data preparation, model development, model evaluation, model post-processing, and model deployment (Baker & Hawn, 2022). However, these decisions may cause a trade-off between having a model that is fair but less accurate or a biased but more accurate model (Aler Tubella et al., 2022). In addition, it is argued that even if developers are aware of the importance of ethics in the development of AI systems, the definition of ethical principles has no impact on the way these developers work (Munn, 2022).

Although the relationship between fairness and bias is not clear in the literature, multiple research papers use the terms interchangeably and attempt separations between technical definitions and moral implications (Baker & Hawn, 2022). This causes clashes between measurements, naming conventions, and world views. The contextual nature of fairness seems to be challenged by the contextual boundaries of machine learning models and their application to a variety of fields. These clashes of fairness notions, world views, consequences, are discussed in the sections below.

2.4.1. Harms and consequences of machine learning models

The literature argues that machine learning models may reinforce discrimination if the datasets reflect societal inequities and historical biases (Salimi et al., 2020). Moreover, multiple papers define their own fairness metrics but lack analysis of their approaches across datasets (Gao & Shah, 2020). Dehouche (2021) argues that pre-trained models may reflect and amplify and reproduce stereotypes by creating correlations between independent labels. The implications of models used in decision-making may produce

unfair outcomes so some authors argue that it is necessary to assess quantitatively the fairness of such decisions (Fletcher et al., 2021).

Some authors argue that the first step to prevent discrimination is to choose a fairness definition in mathematical terms that would allow the measurement of the potential bias (Aler Tubella et al., 2022). The idea behind is that detecting biases would mitigate algorithmic discrimination and build fair models (Giovanola & Tiribelli, 2022). However, biases can come in multiple ways, as they can be completely unknown to the designers or they can be known to exist but not how they manifest or affect the outputs (Baker & Hawn, 2022). Related research highlights the harms that can result from algorithmic bias. For Barocas et al. (2017), harms can be classified as allocative and representational harms. Allocative harms are the result of holding opportunities or resources from specific groups, such as racial bias in ad delivering, gender bias in credit limits, sentencing decisions, identification of patients for medical care, standardised testing. On the other hand, representational harms are the negative systematic representation of a group (Suresh & Guttag, 2021). Such harms include denigration, stereotyping, under-representation, lack of recognition, between others (Baker & Hawn, 2022).

Multiple studies describe how machine learning models can be discriminatory against minority populations. In the United States, medical models were found to discriminate against minorities that couldn't afford medical care (Fletcher et al., 2021). Similar results may happen if an imperfect proxy is used to collect data from populations that present higher cancer rates. The innature identification of higher risk individuals may cause missed diagnoses for underserved communities (Paulus & Kent, 2020). Similar risks are becoming an increasing problem as high-stakes admissions tests are using automated scoring systems (Johnson et al., 2022).

Similar studies describe how bias was found towards male candidates because of the presence of historical discrimination in the training data in the hiring system of Amazon (Salimi et al., 2020). Such a project was abandoned in 2017. Other examples include Tay, the online chatbot of Microsoft, that learned from tweets and started using racist slurs (Chhabra et al., 2021). More recent studies foreseen potential harmful uses, such as misinformation, spam, phishing, fraudulent academic essay, and others for GPT-3, a pre-trained language system developed by OpenAI (Dehouche, 2021).

Probably, the most used example of unfairness in algorithms is the system COMPAS. The system was developed to estimate recidivism rates but was found to predict that black individuals were more likely to commit a crime than white individuals despite sharing similar characteristics in other attributes (Chhabra et al., 2021). The system used data from 7000 arrests and the misprediction led to potentially longer sentences, including those that presented a risk of 0%. Moreover, the algorithm was developed to be race-unaware but other attributes that were correlated with race were included (Paulus & Kent, 2020). The prediction of recidivism rates generated discussions about the fairness of the system. The developers argued that the algorithm satisfied predictive parity, meaning that true-positives and true-negatives were equal for both groups. However, the algorithm did not satisfy equal false-positive and false-negative rates (Brandao et al., 2020). A more comprehensive description of the incompatibility of fairness notions is written in the next section.

2.4.2. Incompatible fairness notions

The amount of fairness notions defined in the literature has led to multiple types of classifications and taxonomies (Ashokan & Haas, 2021; Castelnovo et al., 2022; Green, 2022; Verma & Rubin, 2018). In general terms, algorithmic fairness can be classified in notions based on parity between groups, notions preventing different treatment of individuals, and notions based on causal relations (Baker & Hawn, 2022; Castelnovo et al., 2022). Criteria that is based on causal relations attempts to use domain knowledge to understand the causal structure of an issue (Castelnovo et al., 2022). It can be described as a thought experiment where different scenarios happen where most of the variables are equal but the protected attributes (Piccininni, 2022).

Group fairness aims to reach equality of treatment between groups that can be defined based on their sensitive attributes (Aler Tubella et al., 2022; Castelnovo et al., 2022). Notions for group fairness are derived from Disparate Impact, meaning that no group should be affected by the outcome of a decision-making system (Chhabra et al., 2021). Moreover, group fairness is linked to Demographic parity, Statistical parity, Equalised odds, Calibration and Predictive parity notions (Castelnovo et al., 2022). Some of the flaws of group fairness is that it requires to satisfy conditions only on average of groups, this would lead to unfairness within the groups. Moreover, if two groups behave

differently can it be considered fair to have demographic parity between such groups? (Castelnovo et al., 2022).

Individual fairness is based on the idea that equal individuals should be treated equally despite their differences in protected attributes (Aler Tubella et al., 2022; Castelnovo et al., 2022; Chhabra et al., 2021). Individual fairness is linked to counterfactual fairness, individual calibration, treatment equality, and fairness through awareness. The literature highlights the issue of defining what similar individuals are in the context of individual fairness ((Baker & Hawn, 2022). The issue of defining a similarity metric is considered one of the most challenging aspects of individual fairness (Gupta & Kamble, 2021). Moreover, the idea of treating similar individuals similarly is restrictive as it is not always possible to know what are the best decisions a priori (Gupta & Kamble, 2021).

It has been proved that the satisfaction of more than one fairness notion is not always possible (Baker & Hawn, 2022). For instance, equalised odds and calibration require contradictory requirements to be achieved simultaneously (Piccininni, 2022). Moreover, achieving group fairness, negative class balance, and positive class balances was also proved to be impossible simultaneously (Kleinberg et al., 2016).

The discussion of the impossibility between fairness notions is further expanded with the role of protected attributes in the datasets. While in some cases removing the protected attributes can be beneficial, other cases require such attributes for making fair decisions (Giovanola & Tiribelli, 2022). For instance, if protected attributes are removed, important information may be lost while at the same time, correlations coming from these removed protected attributes may not be fully erased (Baker & Hawn, 2022). Moreover, removing the sensitive features is not sufficient to eliminate bias (Ashokan & Haas, 2021). Additionally, if the protected attributes related to gender are removed, for example, then flipping attributes would not produce counterfactuals (Baker & Hawn, 2022). These decisions are determined on the basis of societal norms and the consideration of what is fair about the issue at stake (Aler Tubella et al., 2022).

2.4.3. The impossibility of algorithmic fairness

The multiple fairness definitions in the literature show that fairness can be a multi-faceted concept (Castelnovo et al., 2022). Some studies argue that having one fairness definition would not be adaptable across contexts since the term may carry different meanings for each different context (Saxena et al., 2020). Moreover, the efforts to achieve fairness may be limited by the “impossibility of fairness” which argues that an algorithm cannot satisfy all desired mathematical notions of fair decision-making (Green, 2022). The definition of fairness is relative and depends on the application and the issue at stake (Dehouche, 2021). For instance, getting a lower recidivism score or a higher credit rating would be beneficial from the perspective of the individual who is being predicted (Paulus & Kent, 2020). The variety of fairness notions may be then a reflection of the impossibility of formalising an unique, universally-accepted definition of fairness (Piccininni, 2022).

No single mathematical definition of algorithmic fairness encapsulates the philosophical views discussed earlier (Green, 2022). Any system will then always encode some belief and make some assumptions about the world as algorithms designers are forced to make assumptions to make real-world decisions (Friedler et al., 2021). Moreover, the lack of a universal definition of fairness cannot be solved either by using human decision-makers to replace algorithms (Paulus & Kent, 2020). If the diverse preferences, social roles, and interests influence how humans conceptualise fairness, then no universally accepted definition will be applicable at all times (Piccininni, 2022). Additionally, if algorithms are trained with data coming from human decisions, it's also not possible to trust the objectiveness of the outputs, meaning that even data-driven decisions may be unfair (Castelnovo et al., 2022). Current formal algorithmic fairness follows a systematic approach to formulate the problem and operationalize fairness around an isolated decision-making process (Green, 2022). Moreover, formal algorithmic fairness focuses on balancing any trade-off between fairness metrics instead of choosing a single one (Green, 2022).

In practice, it seems impossible to be sure if important variables were omitted so statistical methods are used for evaluating fairness of model outputs (Johnson et al., 2022). Formal equality is prone to reproduce the unfairness patterns of the society

(Green, 2022). Because of the diversity of attributes that can be considered and the multiple approaches to model these, there are many different ways to define what is similar (Brandao et al., 2020; Paulus & Kent, 2020). In consequence, the risk of any given outcome becomes model-dependent (Paulus & Kent, 2020). To fix such issues, some authors suggest over-sampling the data from less common groups to prevent any unfair treatment (Baker & Hawn, 2022). However, as previously described, it is not always possible to create fair systems for multiple protected groups with different bias manifestations (Zehlike et al., 2022).

Even if a human decision-maker may have cognitive biases, a biased algorithm may be able to perpetuate discrimination at-scale (M. S. A. Lee & Floridi, 2021). In order to ensure that an algorithm is designed ethically, all development decisions may require to be scrutinised (Aler Tubella et al., 2022). However, a more rigorous definition of unfairness would require understanding the causes for differences in the outputs (Paulus & Kent, 2020). A substantive fairness approach would expect to account for structural inequities, remedies for these, and question if algorithms are an effective tool to make possible change in these relations (Green, 2022).

2.4.4. Escaping the impossibility of fairness

Giovanola & Tiribelli (2022) argue that the philosophical reflection of fairness in the context of theories of justice is rooted in the inability of recognising the quality of value of individuals. To respect an individual it is required to respect their status as an individual and focus on how they can exercise their agency so they can give meaning and purpose to their life (Giovanola & Tiribelli, 2022). If people consider the information about race and gender important for distributive justice, for example, then addressing such historical inequities may require taking into account such attributes in algorithmic fairness (Saxena et al., 2020). To ensure fair equality, eliminating biases may not be enough and taking existing social inequities into account may be required to compensate for them (Giovanola & Tiribelli, 2022).

Multiple authors have defined expected principles to achieve algorithmic fairness as an effort to define what is expected from these mathematical notions. Paulus & Kent (2020) describe principles for accountable algorithms including: responsibility, explainability, accuracy, auditability, and fairness. For a fair representation of an individual, Zehlike et

al. (2020) argue that algorithms should satisfy: individual fairness, group fairness (statistical parity), monotonicity, and utility. Alternatively, Giovanola & Tiribelli (2022) argue that fairness has a distributive and socio-relational dimensions and it has three essential components: fair equality of opportunity, equal right to justification, fair equality of relationship.

Other papers have tried to define solutions that allow the model designer to navigate between multiple fairness notions on demand or to use approaches from other disciplines. Zehlike et al. (2020) define a model with a parameter θ that can be defined consciously and allows the designers to skew the model towards individual or group fairness depending on a specific scenario. A different approach defines a recommendation method based on the Walrasian equilibrium (an economic concept to guarantee the equilibrium of demand and supply) while arguing that it may be possible not to satisfy all criteria of fairness due to contradictions (Xia et al., 2019). Others attempt to define frameworks to define what attributes of a dataset are sensitive and measure their degree of fairness (Y. Li et al., 2021).

As previously discussed, there is not an universally accepted definition of fairness (Zwick, 2019). Despite the numerous papers that research fairness, it is difficult to formulate a working definition of the concept (Brandao et al., 2020). However, the fairness notions defined in a set of short term metrics may cause overall greater harm in the long term, even if they display positive short term outcomes (Card & Smith, 2020).

3. The Wickedness of Algorithmic Fairness

The previous chapter provided a comprehensive literature review on algorithmic fairness, establishing a theoretical framework that encompasses the challenges observed and the contrasting philosophical approaches between positivism and constructivism previously observed in psychometrics. This duality is reflected in machine learning models as well. While multiple authors have proposed multiple fairness notions, more recent authors have argued for a mathematical incompatibility between these metrics and a need to expand the analysis beyond specific decision points. This duality of perspectives may have its roots in the early studies of computer systems, where complex systems and network fairness were defined in terms of resource allocation based on mathematical formulae. In this way, proposing new mathematical notions would not be sufficient to address the potential harms towards groups in disadvantage. Thus, in order to find other approaches to fair algorithmic systems, it is important to acknowledge that the problem is unsolvable in its current form. Building on the findings of the previous chapter, this section explores fairness as a wicked problem through the lens of the wicked problems framework. First, the wicked problems framework is presented to be followed by a discussion of how algorithmic fairness can be considered a wicked problem.

3.1. Wicked Problems Framework

The wicked problems framework emerged when Rittel & Webber (1973) rejected the idea of rationalisation of design. The authors joined a group of dissenters that argued that the design progress is poorly explained in terms of goal settings (Coyne, 2005). Professional work was being criticised because of its way of mimicking the cognitive style of science, which led to it being unsuccessful to work on social problems (Rittel & Webber, 1973). Traditional methods predicted that the best way to work was to follow a top-down process from problem to solution (Conklin, 2006). In this way, the job of a professional was to solve problems that were definable, understandable and consensual (Rittel & Webber, 1973). Social sciences had to rely upon axioms of individualism that led to assume that public welfare comes from the summation of individualistic choices (Rittel & Webber, 1973). However, moving from a rationality based to an empiricist one

shifted the issue from defining a problem rationally to the idea of community consensus (Coyne, 2005).

Rittel & Webber (1973) argue that in a pluralistic society there is not an objective definition of equity and the consequences for it, challenge the tests for efficiency that were used to measure accomplishment. The authors further argue that this led to questioning whether a planning task is the right thing to do. Moreover, as Western societies become more heterogeneous, each group may hold different values making it impossible to define if a group is right or if it should have its ends served. Members of society that don't have their means served may sabotage projects since every stakeholder in a design problem may see their own solution for the problem as correct (Conklin, 2006).

The problems that planners face are different from the ones that scientists deal with, as planning problems are inherently wicked (Rittel & Webber, 1973). While in the linear model of design thinking problems designers first identify the definite conditions in which they can calculate a solution, wicked problems suggest that there are no definitive conditions to these problems (Buchanan, 1992). For instance, designing transportation policies is a wicked problem because they are loosely formulated, they are later redefined, they depend on the point of view of the designer, and they may have multiple solutions. In order to define wicked problems, Rittel & Webber (1973) list the following ten properties:

1. There is no definitive formulation of a wicked problem: There may be an exhaustive list of all possible solutions. However, the problem cannot be defined until the solution is found.
2. No stopping rule: Without a criteria to define when there is sufficient understanding of the problem, there can always be opportunities to try further options.
3. Solutions are not true or false, they are good or bad: The judgement of the solution will differ between groups and stakeholders depending on their value sets. Solution assessment is often described as better or worse.
4. No ultimate test of a solution: Any solution will have its own consequences over a period of time.

5. Every solution is a one-shot operation: Every solution counts as its consequences may have implications that cannot be undone.
6. Not enumerable set of solutions: There is no criteria to prove that all solutions have been defined.
7. Every wicked problem is unique: Every new problem has its own set of distinguished properties that may not be applicable to previous ones.
8. A wicked problem can be a symptom to another problem: In the search of causal explanations, the removal of a cause may pose another problem where the previous one is a symptom.
9. The choice of the explanation determines the nature of the problem resolution: Without having a rule to define the correct explanation, every person chooses the most plausible explanation to them.
10. No right to be wrong: The planner has liability for the consequences on all the people affected by the actions taken to solve the problem.

For Conklin (2006), the natural forces that make collaboration difficult in which knowledge is scattered, are the forces of fragmentation. Problem wickedness is a force of fragmentation that most projects have. If designers fail at recognising the wickedness of a problem, inappropriate methods and tools may be applied to them. As a consequence, in order to cope with wicked problems, the two most common mechanisms are to study the problem or to tame it. However, the author argues that tame problems are not necessarily technically simple. On the other hand, social complexity is another force of fragmentation where the more parties get involved in collaboration, the more socially complex it becomes. Social complexity is then a function of the structural relationship between stakeholders where the rule of everyone thinking and acting the same doesn't apply. In that way, social complexity makes wicked problems even more wicked and the issue is then the failure of realising the dynamics of the systemic mess. In other words, a project is about reconciling what is needed and what can be built.

Coyne (2005) criticises the wicked problem theory by arguing that all problems can be considered wicked problems. A tamed problem is only a microworld that designers can use through interpretive skills, that derives from the level of abstraction in which a problem is formulated. However, despite the critics other authors suggest that taking a

systemic view of the problem where the whole group gains certain literacy in the language of coherence, a shared understanding and commitment to the problem can be created (Conklin, 2006).

3.2. Algorithmic Fairness is a Wicked Problem

As discussed in previous chapters, the vast literature on algorithmic fairness has presented the topic as contested and difficult to achieve. In order to find other approaches to achieve algorithmic fairness, it is important to acknowledge that the problem is unsolvable in its current form. The present section analyses algorithmic fairness from the perspective of wicked problems and creates a starting point for a discussion beyond the idea of fairness notions.

1. There is no definitive formulation of a wicked problem

Rittel & Webber (1973) explain that it is possible to provide an exhaustive formulation with all necessary information required by the problem-solver for tackling tame problems. Unlike this, most of the solutions used to achieve algorithmic fairness have not been able to find a widely accepted definition of fairness. In psychometrics, it was already argued that research was not able to provide a definition for fairness to indicate the existence of unfairness (Cole & Zieky, 2001). In complex systems, the intuitive definition used in most solutions was considered to be fair for some states and not others (Queille & Sifakis, 1983). In networks, the maximin definition would only be fair if a network is seen as a whole (Chan & Zukerman, 2002).

Wicked problems require having an extensive list of all possible solutions to anticipate all possible questions (Rittel & Webber, 1973). The lack of a universally accepted definition of fairness has led to the development of numerous models, metrics and concepts of algorithmic fairness (Ashokan & Haas, 2021; Castelnovo et al., 2022; Green, 2022; Mehrabi et al., 2021; Verma & Rubin, 2018). These definitions build upon one another in order to include all possible scenarios that previous notions did not consider. However, algorithmic fairness also relies on the values and beliefs that the developers of algorithms encode in their systems (Friedler et al., 2021). Moreover, there are many ways to define what is similar or what protected attributes to consider as they may rely on cultural differences and become model dependent (Paulus & Kent, 2020).

If we trace the problem of achieving algorithmic fairness to the original problem of defining what fairness is, then probably the solution of algorithmic fairness would be found. However, as previously mentioned, multiple authors have tried to define the term and no consensus has been achieved. Current fairness notions rely on mathematical equations to measure fairness and attempt to convince other approaches to be less effective. Such an approach would not be possible while trying to achieve fair treatment with all stakeholders included in the use of such systems. It wouldn't be considered a solution to tame the problem as it is, since it would require creating a more specific setting. Such a setting would then leave out certain important considerations to take while working with algorithms that are available across cultures. The formulation of algorithmic fairness in its current form looks then, wicked in nature.

2. There is no stopping rule

As mentioned before, there are multiple fairness notions mentioned in the literature of algorithmic fairness. According to Rittel & Webber (1973), there is no criteria to know if sufficient understanding of a problem has been acquired. Any kind of additional investment into the topic would lead to developing a fairness notion that benefits a bit better to the majority of the group. For instance, in some domains, including algorithmic fairness, it was shown that late approaches tried to achieve dynamically more than one fairness notion at once by adjusting variables. However, it is also argued that these fairness notions may not be compatible across datasets or contexts (Gao & Shah, 2020; Saxena et al., 2020). This allows multiple questions to ask when a solution is good enough to solve the problem. After all, as argued before, choosing a definition of fairness would rely on subjective criteria (Apt et al., 1988).

3. Solutions are not true or false, they are good or bad

It was shown that in other domains the output of a process can be seen positively or negatively depending on the values of different stakeholders (Bisio & Marchese, 2014; Cole & Zieky, 2001; Kwiatkowska, 1989). In psychometrics when discussing the Cleary model and further contested models, every solution seemed better than the previous one. The intuitive definition of complex systems that deal with infinite processes was questioned on its means of application and how it was open to interpretation. The discussion of Maximin fairness in networks reflected a similar issue when every new

technology added a layer of complexity to the problem by creating more distributed systems that needed a better bandwidth allocation. Such a discussion is also presented in wicked problems, where multiple stakeholders with equal knowledge are entitled to judge the solutions, while no one has the power to define the correctness of each (Rittel & Webber, 1973).

This problem is enhanced by the faulty assumption that every individual would share the same idea of social justice (Zollers et al., 2000). The literature shows that these algorithmic systems will encode some beliefs of the world (Friedler et al., 2021). Moreover, it is not possible to trust its objectiveness if they are trained in human data (Castelnovo et al., 2022). In this way, even the solutions considered for algorithmic fairness are based on societal norms about what is fair of the issue at stake (Aler Tubella et al., 2022). This means that not all fairness definitions will be adaptable across contexts and datasets (Gao & Shah, 2020; Saxena et al., 2020). As highlighted in the wicked problems framework, any solution of a wicked problem would only be assessed as better or worse than other ones. This is a clear issue that is part of the algorithmic fairness discussion and even if some authors may claim that their solution should be used, they are still encoding their own values. This may be a reflection of social complexity as a force of fragmentation discussed by Conklin (2006). As machine learning models are used in more domains, more stakeholders get involved and thus, more socially complex the problem of algorithmic fairness becomes.

4. No ultimate test of a solution

The solution of tame problems are easy to determine, in a way that they are under control of the people involved in the problem (Rittel & Webber, 1973). However, as seen in psychometrics, some individuals may rely on their own personal expectations or their own interpretations of justice principles to define whether a solution is fair or not (Crocker, 2003). Moreover, as argued in the previous section, fairness definitions may not be adaptable across contexts or datasets (Gao & Shah, 2020; Saxena et al., 2020). Algorithmic fairness deals with a more difficult problem when it is argued that any implementation may be limited by the impossibility of fairness, meaning that any algorithm cannot satisfy all mathematical notions at once for fair decision-making (Green, 2022). What would be the ultimate test of an algorithmic fairness solution if the

existing notions are not mathematically compatible? This question is probably part of the formulation problem mentioned earlier.

5. Every solution is a one-shot operation

For Rittel & Webber (1973), every solution for wicked problems has its own consequences where every outcome cannot be undone or may leave traces. Such a problem in algorithmic fairness can be seen in the literature too. For instance, the system COMPAS, which was found to predict recidivism rates favouring white individuals over black ones (Chhabra et al., 2021). Another example of such traces that an implementation may leave are the hiring systems by Amazon (Salimi et al., 2020), the chatbot Tay of Microsoft (Chhabra et al., 2021), and the potential misuse of GPT models (Dehouche, 2021). The harmful consequences of these projects were discovered and studied after their implementation. This is in line with the idea that machine learning models encode the values of their developers and that algorithm designers may be forced to make assumptions to make real-world decisions (Friedler et al., 2021).

The argument here is that these algorithms included their own perceptions of the world in their implementation. This means that even without noticing, they may cause harmful consequences in the people who use them or rely on their decisions. As the wicked problem framework argues, every trial counts and every attempt to reverse such decisions pose a new set of wicked problems (Rittel & Webber, 1973).

6. Not enumerable set of solutions

As previously discussed, algorithmic fairness has no definite set of solutions. Every new fairness notion would make things a bit better than the previous one. As mentioned by Rittel & Webber (1973), the problem-solver may arrive at a description that requires that two incompatible views happen at the same time. That is the case of the validity models in psychometrics. The Cleary model would not be fair according to the Darlington model and vice versa (Cole & Zieky, 2001). In a similar note, algorithmic fairness notions have been demonstrated not to be all compatible mathematically with each other (Baker & Hawn, 2022; Kleinberg et al., 2016). In addition, several authors have created exhaustive lists of notions of fairness, varying in the notions under consideration (Ashokan & Haas, 2021; Castelnovo et al., 2022; Green, 2022; Mehrabi et al., 2021; Verma & Rubin, 2018).

7. Every wicked problem is unique

The biases of algorithmic systems may come in multiple ways, they can be either unknown to the designers or known to exist while not knowing how they manifest (Baker & Hawn, 2022). The fairness notions developed for psychometrics may not be directly applicable to algorithms and they may not include the complexities presented by social interactions they faced. What has been found to be similar to other domains was the initial development of fairness notions that followed an intuitive mathematical approach. These initial definitions were later followed by notions that were incompatible with the initial definition. After the incompatibility between fairness notions was acknowledged, multiple authors tried to allow dynamic adaptation to every notion depending on the context. Once this seemed as an impossible task that would not fix the problem, research turned into understanding the diversity of experiences by the people who were affected by these fairness notions. In this way, the complexity of algorithmic fairness would be ill-advised to follow the solutions that were implemented in other domains as highlighted with a planning scenario by Rittel & Webber (1973).

8. A wicked problem can be a symptom to another problem

As seen in psychometrics, fairness is a contested topic with multiple definitions and notions that are not always compatible with each other. Some papers have tried to define solutions that adapt dynamically to the required fairness notions depending on the problem at stake (Zehlike et al., 2020). Maybe the discussion on the impossibility of algorithmic fairness is rooted into the impossibility of formalising a universally-accepted definition of fairness (Piccininni, 2022). This would mean to discuss the problem in a higher level which, as argued by Rittel & Webber (1973), would increase the difficulty of the issue. As discussed in the implications of the development of selection tests, fairness may be considered an individual interpretation of certain principles of justice that the individual holds. This leads to wonder what theories of fairness exist and how they are related to the existing fairness notions. A theoretical framework of fairness would help to understand how it is connected with principles of justice and probably open new approaches to try to escape the wickedness of algorithmic fairness.

9. The choice of the explanation determines the nature of the problem resolution

It has been seen that previous research in other domains has gone through a similar process where initial fairness notions were dominant and preferred to later be challenged by solutions following different considerations. As Rittel & Webber (1973) argue, every explanation for a wicked problem may be picked in order to fit best the intentions and action-prospects of the problem-solver. Every fairness model used in psychometrics would solve the biases that the authors found in previous models (Baharloo, 2013; Berry, 2008). Such a debate intensified when it was proved that some tests could be valid for everyone but unfair to certain groups depending on what definition of fairness was considered (Cole & Zieky, 2001). If the diversity of preferences, social roles, and interests influence how an individual conceptualises fairness then no definition can be applied at all times (Piccininni, 2022).

Other researchers argued that it is faulty to assume that everyone shares the same idea of social justice (Zollers et al., 2000). Multiple fairness notions are proposed while being in constraint of the definitions of other papers (Bisio & Marchese, 2014). For instance, in the literature, the existing fairness notions are often categorised in different ways. Individual fairness metrics may not be compatible with group fairness metrics, while causal notions may have a different way of application. At the same time, it was seen that the term fairness was used interchangeably with bias (Baker & Hawn, 2022). The analysis and world view that every author used to define these fairness notions will be the determining factor to try to solve the wicked problem of algorithmic fairness. If there are so many contested notions, escaping the wickedness of algorithmic fairness would go beyond the definition of a new metric that attempts to achieve all the others.

10. No right to be wrong

As shown in this whole discussion, every new fairness notion forces the designers to make assumptions of the real world. Every solution developed may have unintended consequences that would have a high impact by reproducing the unfairness patterns in society (Green, 2022). For instance, a biased algorithmic system may perpetuate discrimination at-scale (M. S. A. Lee & Floridi, 2021). Moreover, the effects of these systems matter the most to the people that are touched by these actions (Rittel & Webber, 1973). As discussed in the previous sections, algorithmic systems may have different kinds of harms that may cause negative consequences in disadvantaged

groups. This problem raises questions about to what extent unfairness should be tolerated, who may make such decisions, and what approaches can be followed to achieve fairness in algorithmic systems that are applied in multiple domains and fields.

The consulted literature presents all previous arguments to consider algorithmic fairness a wicked problem. A next potential step would then be to wonder if fairness itself is a wicked problem, however, as Rittel & Webber (1973) argue, a wicked problem can be considered to be a symptom of another problem that could also turn out to be wicked. Instead, some authors that have suggested that algorithmic fairness is a wicked problem argue for shifting focus from solutions to processes (Scantamburlo, 2021). This would mean to place more importance into listening to the needs of all stakeholders rather than focusing on the outcomes. In this way, to address the complex issues of algorithmic fairness, stakeholders should not work in isolation and instead, a strong participation from multiple players would be enabled by sharing a robust understanding of the goals to achieve (Woodruff et al., 2018).

The described fairness notions may cause greater harm in the long term if they are based on short term metrics (Card & Smith, 2020). Such a problem is enhanced as a biased model may perpetuate discrimination at-scale (M. S. A. Lee & Floridi, 2021). For instance, it was shown that language models may reflect and amplify stereotypes and that they can be used for harmful uses (Dehouche, 2021; Kirk et al., 2021; Ray, 2023). If biases come in multiple ways that cannot be completely known to the designers then eliminating them would not be enough to ensure fair equality (Baker & Hawn, 2022; Giovanola & Tiribelli, 2022). This situation is a clear indication that fairness metrics may not be enough to measure algorithmic fairness and eliminating biases may not necessarily solve this problem either.

As other authors have previously argued, the mathematical formulations of fairness discussed in the literature overlook the contextual, ethical, and philosophical meaning of fairness (Card & Smith, 2020; Green, 2022). A different approach than the previously presented would be to judge a model against ethical and legal principles to satisfy fairness criteria (Fletcher et al., 2021). However, it is important to consider the relative nature of fairness as a concept and the difficulty to formulate a working definition of the concept (Brandao et al., 2020; Dehouche, 2021). Scantamburlo (2021) suggests that if

algorithmic fairness solutions are inherently imperfect, then a shift in focus to processes rather than solutions is needed to enable collaboration between different stakeholders and minority representatives. Such a collaborative approach would then need to consider the following critical issues:

- There is no widely accepted definition of fairness that can be applied to multiple contexts. Fairness is influenced by the context and the level of analysis.
- The lack of an accepted definition has led to the development of multiple fairness notions that are mathematically incompatible with each other. These fairness notions may also be incompatible across datasets.
- The existing mathematical fairness notions lack an ethical and philosophical foundation. They are unable to account for the systematic and historical patterns of discrimination present in the data.
- The definition of ethical principles alone may have no impact in the way the developers of algorithmic systems work.
- The mathematical fairness notions focus on single decision points and overlook the systematic issues of society.
- Models reflect the stereotypes, historical biases, and patterns of discrimination from their developers and data.
- Models force developers to make assumptions of the real-world for decision-making.
- It is not possible to trust the objectiveness of these models as they are trained with human data.

4. Theoretical Framework: Shifting Focus

The previous chapter presented an in-depth analysis of algorithmic fairness and how it can be considered a wicked problem in its current form. The chapter concluded with a list of critical issues that need to be considered in order to develop a collaborative framework for the implementation of fair algorithmic systems. In this way, in order to overcome the issues caused by the lack of a widely accepted definition of fairness, it becomes crucial to move away from the current mathematical notions of fairness towards a philosophical understanding of the concept. However, as will be explained in the following sections, fairness is not the default term in the literature of philosophy, and instead several authors have discussed *justice* as an appropriate term for the development of theories. As a result, this chapter introduces two of the most prominent modern theories of justice.

The remaining issues listed in the previous chapter require a broader level of analysis. This means to move from immediate outcomes to instead consider the systematic problems in society that cause social inequalities. Substantive algorithmic fairness provides a framework that can be helpful in taking these considerations into account to provide an agenda for algorithmic justice. Substantive algorithmic fairness considers the disparities grounded in social hierarchies and the restricted benefits for the individuals judged negatively. It aims at promoting justice in practice by identifying social inequalities, defining potential reforms, and understanding the role of algorithms to support these reforms. By including a thorough description of the theories of justice and the Substantive algorithmic fairness framework, this section aims to provide a solid theoretical foundation that would contribute to the discussion of the development of fair algorithmic systems, away from the wickedness of algorithmic fairness.

For the convenience of the reader, Table 1 provides a brief summary of the definitions used in the following sections. This summary provides a clear distinction between justice, fairness, formal equality, substantive equality and bias, setting the stage for a comprehensive description of the theories and frameworks presented. Through this approach, the author aims to create an insightful narrative that guides the reader to a deeper understanding of the key concepts involved in the journey towards the development of fair algorithmic systems.

Term	Definition
Justice	Whether an entity adheres to rules, standards, or laws and denotes a conduct that is morally required (Goldman & Cropanzano, 2015).
Fairness	The way an individual responds to the perception of justice principles and the evaluative judgement given to them (Goldman & Cropanzano, 2015).
Formal Equality	The equal treatment for individuals based on their attributes or behaviour at a particular decision point (Green, 2022).
Substantive Equality	The identification and remediation of social hierarchies that generate disparities in social and material resources (Green, 2022).
Bias	The outcomes that are less favourable to individuals within groups where groups have no relevant differences that would justify such harms (N. T. Lee et al., 2019).

Table 1. List of definitions.

4.1. Theories of Justice

Even though the terms *justice* and *fairness* are used interchangeably in research, fairness is not the standard term used in philosophy (van Nood & Yeomans, 2021). For Goldman & Cropanzano (2015), it is important to make a distinction between both terms. The authors argue that *justice* refers to whether an entity adheres to rules, standards, or laws and denotes a conduct that is morally required. Whereas, *fairness* refers to how individuals respond to the perception of these rules and the evaluative judgement given to them. Moreover, as multiple languages do not have the distinction between both terms, multiple scholars discuss if the idea of *Justice as fairness* by John Rawls can be considered as a basis for a theory of justice (Esmer, 2021).

Esmer (2021) gives a historical view of the development of justice. According to the author, there were two approaches to the theory of justice during the Enlightenment period. The first one paid attention to creating just institutions and held the idea of social contract. The main representatives of this approach are Thomas Hobbes, John Locke, Jean-Jaques Rousseau, and Immanuel Kant. John Rawls and his theory of justice is influenced by this view. On the other hand, John Stuart Mill, Karl Marx, and Jeremy Bentham aim to eliminate injustices in the world by focusing on social realisations and interactions between people. This alternative approach has influenced the *Capabilities*

approach proposed by Amartya Sen. In the following sections, both approaches will be explored.

4.1.1. Justice as fairness

The publishing of *A Theory of Justice* by John Rawls caused the end of the dominance of utilitarianism on political philosophy as a way to focus on welfare maximisation as the main objective of policies (Esmer, 2021). With his theory, Rawls aimed at moving the social contract theory of Jean-Jacques Rousseau, John Locke, and Immanuel Kant to a higher level of abstraction and generalisation (Esmer, 2021). A social contract theory focuses on what political principles would be unanimously agreed to be respected by individuals in a society (Robeyns, 2009). Justice as fairness explores how to ensure fairness between individuals in a given society and how to determine just institutions with a framework to define principles of justice (Esmer, 2021). In this theory, the principles of justice are given a higher priority over a single principle for all virtues (Das, 2021).

For Rawls, a society is an association where in the relations of individuals, most of them comply with a set of rules or principles (Rawls, 1971). The benefits and burdens coming up from this social cooperation arise questions of justice (Das, 2021). So for Rawls, a well-oriented society is one that has been “designed to advance the good of its members and effectively regulated by a public conception of justice” (Rawls, 1978). Justice is then defined by “the role of its principles in assigning rights and duties and in defining the appropriate division of social advantages” (Rawls, 1996).

In Justice as Fairness, justice is a political concept that comes from the principles of distributive justice determined by the members of society and it is provided to the extent that the members of society adhere to these fundamental principles (Esmer, 2021). Moreover, in order to develop a Kantian conception of justice, Rawls suggests that the doctrine of Kant must be given a procedural interpretation by the definition of the original position (Rawls, 1978). This original position constitutes the principles of justice that form the structure of society and are chosen by equal people on equal terms in order to balance the advantages, fundamental rights, duties and benefits of the members of the community (Rawls, 1971).

The original position can be considered the appropriate status quo in which fair agreements are made (Das, 2021). So then, for Rawls, the correct way to determine which goods are apportioned and in which ways is by making such decisions behind a veil of ignorance (van Nood & Yeomans, 2021). The veil of ignorance, according to Rawls (1971), allows only limited knowledge to individuals about all the variations of social circumstances that would have an influence in their rights and treatment in society. This leads to having choices that are reasonable for oneself and others. As a consequence, the defined principles under the veil of ignorance ensure that there are no advantaged and disadvantaged individuals as everyone is equally situated and no one can design the principles to favour their own condition (Rawls, 1971). However, the parties involved in the definition of principles know the general facts of society, human psychology, and the relation between individuals and their social background (Robeyns, 2009). In Justice as Fairness, individuals are driven by the interest of having the capacity for a sense of justice and to pursue a conception of the good (Esmer, 2021).

The principles that these free and rational individuals would accept in an initial position would regulate all further agreements (Rawls, 1971). By defining the principles of justice under a veil of ignorance, these principles become universal and the possibility of contingency is ruled out (Das, 2021). This means that the demand of justice as fairness is at par with the demand of impartiality in the original position (Das, 2021). However, Rawls argues that it is important to consider that the original position is not an actual historical state of affairs and instead should be considered a hypothetical situation that leads to a conception of justice (Rawls, 1971).

Rawls suggests two principles that are likely to be agreed in the original position. The first principle requires equality when assigning basic rights and duties (Rawls, 1971). This principle relies on the social primary goods approach, which are goods that anyone wants regardless of what others receive (Das, 2021). The individuals that can be considered the least advantage are the ones that have the lowest index of primary goods (Das, 2021). The distributed primary goods that Rawls propose can be classified in the following groups: 1) basic rights and liberties, 2) freedom of choice of occupation, 3) the rights of holding a position of responsibility, 4) income and wealth, 5) social bases of self-respect (Robeyns, 2009). Rawls argues that the proposed social goods would be preferred by all parties despite the differences between their life plans (Esmer, 2021).

The second principle argues that social inequalities are only just if they compensate everyone, especially the ones who are the least advantaged members of the society (Rawls, 1971). On the application of the second principle, Rawls uses the economic term of Pareto optimality which refers to an optimal arrangement where it is not possible to make a group of people better off without making another better off (Esmer, 2021). In that way, these inequalities need to be arranged so that they are expected to be to everyone's advantage and attached to positions that are open to everybody (Rawls, 1971). The two principles resulting from the original position have been interpreted as the maximin welfare function (Hall, 2021). In the maximin rule, the alternative of the worst outcome is taken which is superior to the worst outcome of others (Rawls, 1971).

The aim of Justice as fairness is not to create a new form of government but instead to choose moral principles (Esmer, 2021). This concept of moral agreement is based on the Kantian concept of moral universalizability (Das, 2021). By following both principles (justice and fair equality of opportunity), everyone is provided with the same understanding of self-respect (Robeyns, 2009). Thereby, the fundamental subject of justice is the way in which the social institutions distribute the fundamental goods and determine the division of social advantages (Esmer, 2021). This leads to assuming as reasonable that no individual should be advantaged or disadvantaged because of natural fortune or the choosing of social principles (Rawls, 1971). This view conceives justices as *justitium*, a form of deontological justice that takes rules and principles as the base of the theory (Das, 2021). However, considering that the reliance on formal rules is a representation of formal justice, argues that there may not be impartial institutions.

4.1.2. The capabilities approach

Amartya Sen develops his idea of justice as *justitia*, with which he favours some form of utilitarianism and denies the deontological approach of Justice as fairness (Das, 2021). This approach is mostly focused on enhancing justice based on how people live in the society and not only on stringent rules (Das, 2021). It differs from Justice as fairness as it doesn't aim to achieve a perfect fair society and instead it aims at reducing injustices (Esmer, 2021). The approach of Sen makes interpersonal comparison by focusing on the opportunities of individuals and who they want to be instead of the primary goods they receive (Robeyns, 2009). For Sen, functionings are achievements and capability is the

ability to achieve these (Sen, 1988). Being healthy, sheltered, not mentally ill, engaged in social relations, able to achieve a work-life balance are some examples of functionings (Robeyns, 2009). In this way, justice should concern the capabilities instead of functionings if individuals have similar opportunities and can decide on their own to take responsibility for them (Esmer, 2021).

The subject of the capability approach does not limit to the structure of society as it includes primary goods, social norms, interpersonal relations, and common attitudes (Esmer, 2021). For Sen, the primary goods approach would be insufficient to consider the interpersonal differences between individuals while converting goods into functionings (Robeyns, 2009). Sen (1980) argues that as there is evidence that the conversion of goods into capabilities varies between individuals, then the equality of this conversion is far from being an equality resulting from having individuals that are very like each other. The capability approach is then about freedom and the creation of an environment where individuals can flourish, in other words, about developing the capabilities people need to follow the life they value (Walker, 2005). In this way, individual claims of justice are to be assessed based on the freedoms they enjoy to choose between the possible ways of living they may value (Sen, 1990).

For Sen, the principles of justice should be defined in terms of the lives and freedoms of the individuals that are part of the society instead of the institutions (Das, 2021). The capability approach is then a “natural extension of Rawls’s concern with primary goods shifting attention from goods to what goods do to human beings” (Sen, 1980). Thus, it is more appropriate to focus on the measurement of capabilities and freedoms along with the primary goods (Das, 2021). Sen may be considered a contemporary supporter of Aristotle by arguing that material goods alone are not important and can only be judged in relation to human functions (Esmer, 2021). For Sen, capabilities cannot be bad in themselves and can only be judged as such based on their use (Saito, 2003). What matters is then to examine what was the reasoning for demanding justice as multiple reasoned positions may exist (Sen, 2008).

The difficulty in Rawls’ approach would be to realise if the plurality of reasons of justice would allow a set of principles to emerge from the original position (Das, 2021). The most important critique from Sen towards Justice as fairness refers to how the approach

of primary goods neglects human diversity (Esmer, 2021). For Sen, a workable solution doesn't require social unanimity but instead crucial attention to public discussion where every individual is seen as a participant in change (Walker, 2005). In the assessment of social issues, multiple critical functionings and values could be taken into account so it may be impossible to reduce them into a commensurable magnitude (Esmer, 2021). As a consequence, Robeyns (2009) suggests that in order to assess the capabilities of an individual we could: 1) analyse the capability inputs, including the system where the individual functions, the interaction with others, and the primary goods; 2) scrutinise the social, environmental and personal conversions of these primary goods into capabilities; 3) if any conversion factor lowers income, for example, then this would be an argument to provide extra resources or social policies; finally 4) scrutinise the social constraints of the choices that individuals make, this could be done later with the evaluation of the distribution of resources and policies. In this way, the functionings of an individual can be shaped by relative advantages in society and enhanced by policy environments (Walker, 2005).

For Sen, the contractarian approach of Kant and Rawls denies the role of emotions and follows reason alone in the search for justice (Das, 2021). Moreover, for Sen, the social contract theory of justice must allow incompleteness in a social arrangement and should not be considered an universal idea as assumed by Rawls (Das, 2021). As a consequence, Sen leaves his approach open and vague so communities can decide what capabilities they consider valuable (Walker, 2005). Esmer (2021) argues that Sen follows a particularist approach which regards justice as incommensurable because of its plural nature and holds each particular understanding as sovereign. In contrast, Rawls follows an universal idea of justice that is not influenced by human experiences, an approach also known as universalist approach. As a consequence, Sen criticises the universalist approach in two ways: 1) universalists fail to acknowledge the presence of power in human relationships and 2) they prioritise the idea of a perfect justice assuming that just practices will follow (Esmer, 2021). However, Rawls also discussed the scope of Justice as fairness to only be limited to constitutional democracies (Robeyns, 2009).

The term *justitium* focuses on abstract rules that are applied to an ideal situation where emotions and shared beliefs are ignored (Das, 2021). For example, if utility is considered to evaluate justice, the focus is on the mental reaction rather than the person's

capabilities (Sen, 1980). On the contrary, justitia, considers comparative arrangements, shared beliefs, emotions, and reason as a way to determine forms of life (Das, 2021). This shift from justitium to justitia resulting from the capabilities approach would mean a shift from deontology to utilitarianism (Das, 2021). However, despite denying Kantian deontology by denying Rawls approach, Sen, like Kant, admits individual freedom by arguing that the capacity of choosing is a significant human aspect of life.

4.2. Substantive Algorithmic Fairness

Green (2022) proposes the term *Substantive algorithmic fairness* as a new direction to use algorithms to promote justice in practice, away from formal mathematical models of fair decision-making. The author argues that algorithmic fairness mirrors the problematic tendencies of anti discrimination as a mechanism to achieve equality. These tendencies include focusing on bad actors, individual axes of disadvantage, and a limited set of goods. In this way, algorithms that only satisfy fairness standards end up exacerbating unjust institutions.

Green (2022) discusses a central tension in egalitarian theory, the tension between formal and substantive equality. Formal equality asserts that if two individuals have the same status in one respect, they should be treated equally in that respect. This view is limited, according to the author, as this view restricts the analysis to a specific decision point and cannot account for the inequalities that surround them. This leads to the argument that formal equality is prone to reproduce existing patterns of injustice. In algorithmic systems, formal fairness is defined as a technical attribute of the algorithms that formulates the problem around the input and outputs of a specific decision point. The author then describes two responses to the impossibility of algorithmic fairness. The fair contest response argues to measure fairness solely on the likelihood to exhibit an outcome of interest. Therefore, in order to advance algorithmic fairness, the accuracy of these systems should be increased. This view does not take into account the group differences product of oppression. The second response is the formalism response. This view focuses on balancing the tradeoffs of the multiple existing fairness metrics and disregards structural reforms related to racial disparities. This view provides clarity in a limited scope of analysis. What it seems to be impossible then, according to the author, is to satisfy all fairness notions in an unequal society.

Substantive equality, on the other hand, focuses on identifying and redressing social hierarchies that produce inequalities in social and material resources (MacKinnon, 2011). For Green (2022), having social hierarchies at the heart of the inequality problem, neither treating everyone equally nor giving special treatment to the disadvantaged will lead to greater equality. In the context of algorithms, the aim is not to incorporate them into a formal mathematical model and instead include the relational and structural considerations of a decision point into the analysis. For the author, substantive algorithmic fairness should propose reforms to relational and structural inequalities. A relational response would reduce the disparities coming from social hierarchies. This way, the dilemma of treating everyone equally or not would be avoided by not translating the differences between people into disparities in normative significant attributes. A structural response would reduce the scope on which decisions act over social disparities. With this response, the dilemma of treating everyone equally or not would be avoided by preventing decision-making structures to harm individuals who exhibit attributes deemed as negative.

Green (2022) proposes a complete approach to promote equality in decision-making processes where discrimination or inequality are concerned. This approach includes a three-step strategy to promote equality. The first step is to identify the conditions of hierarchy and how social arrangements reinforce those conditions (MacKinnon, 2011). If there are no disparities in the data and if decisions do not exacerbate social hierarchies, then formal algorithmic fairness is the appropriate path. Otherwise, reforms based on formal equality will be insufficient. After diagnosing inequalities, the second step considers what reforms can remediate these. The first approach follows the relational response, to reduce disparities that reflect social hierarchies fed to the decision-making process. The second approach follows the structural response, to reduce the extent in which decisions exacerbate social hierarchies. The final step is to analyse if algorithms can enhance or facilitate the identified reforms. While decentering technology from the discussion of injustice, this analysis will reveal new roles for algorithms to combat oppression. Figure 1 shows the original flow chart for substantive algorithmic fairness.

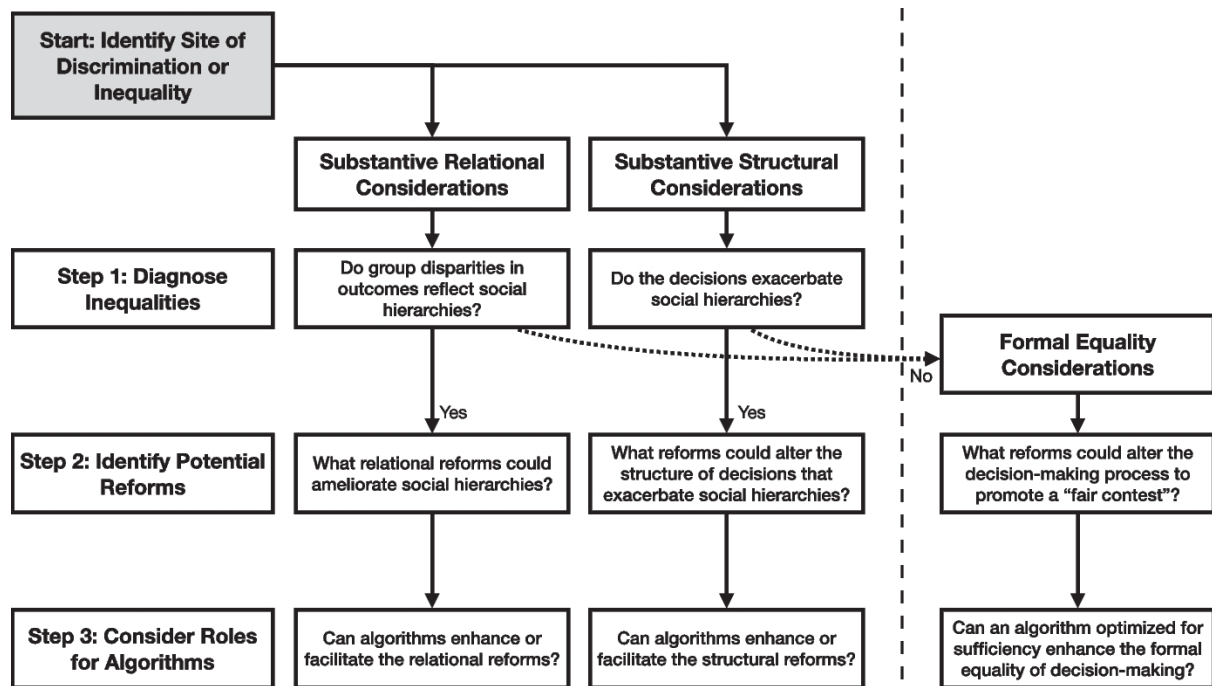


Figure 1. Substantive Algorithmic Fairness framework (Green, 2022).

For Green (2022), to promote justice in practice is necessary to account for relational and structural inequities, theories to remedy those inequities, and consider the role of algorithms to support such reforms. Substantive algorithmic fairness requires a wider scope where disadvantaged communities have a voice to answer the questions diagnosing inequalities. This approach offers a method to develop incremental reforms to push policy to achieve greater substantive equality. This includes considering that algorithms may not be productive tools for promoting certain reforms. In this way, the impossibility of fairness resulting from formal fairness may be avoided.

5. Analytical Framework: A Proposal for Algorithmic Justice

This chapter introduces an analytical framework that proposes a novel perspective on algorithmic justice. Building on Green's framework of Substantive algorithmic fairness (SAF), this chapter uses the presented theories of justice to further enhance the understanding and application of justice in algorithmic systems. In doing so, this chapter aims to provide a novel approach to address the critical issues that current implementations of formal algorithmic fairness present and are discussed in Chapter 3. The following sections explore the theoretical foundations, the integration of justice theories, and the implications for algorithmic justice, providing a robust analytical framework that can represent a new direction into the discussion of fair algorithmic systems.

The Substantive algorithmic fairness (Green, 2022) aims to reform the methodology of algorithmic fairness away from mathematical models, in order to promote justice through substantive equality. This demand for justice, represents a need for those impacted by the outputs of AI systems and a source of philosophical insight (Gabriel, 2022). For Munn (2022) AI justice provides a useful term to expand the ethical scope of intervention. AI justice, reframes the discussion around AI ethics to consider that the moral properties of algorithms are not internal to the models and instead are a product of the social systems where they operate (Gabriel, 2022). A just society will attempt to eliminate the impact of unchosen features on their life prospects (Gabriel, 2022). In this way, SAF aims to include the relational and structural considerations of a decision point into the analysis and move beyond the formal implementation of mathematical algorithmic notions (Green, 2022). This broader analysis may allow a more deep reflection of the notion of *human* and what it means to benefit humanity, by considering the historical issues of disadvantage towards certain groups deemed as less human (Munn, 2022).

In practice, AI justice would mean to consider the views of groups that are affected by AI but are not normally consulted (Munn, 2022). The application of philosophical principles would allow moving towards considering how AI tools can mitigate the effect

of biases that exist at the social level (Gabriel, 2022). A measure that is sensitive to heterogeneous political and social conditions would better address representational harms (Lundgard, 2020). The relational response of Substantive algorithmic fairness would mitigate the extent to which oppressed groups exhibit attributes that are deemed as negative (Green, 2022). This response would not provide differentiated treatment to individuals and instead focus on reducing the disparities grounded in social hierarchy. By addressing these disparities, the representational harms in words of Crawford (2017), the allocative harms that include the distribution of resources and goods, may also be addressed.

This thesis builds over the Substantive algorithmic fairness framework introduced by Green (2022). As the author suggests, this framework suggests only a sequence of questions to promote reforms of algorithmic justice. As a consequence, this proposed novel approach consists of an early attempt to apply the framework in machine learning models. Probably, the main difference to the original proposal is the inclusion of theories of justice in the whole framework. The following figure shows how all different theories may be connected to create a framework to implement algorithmic justice in transformer models.

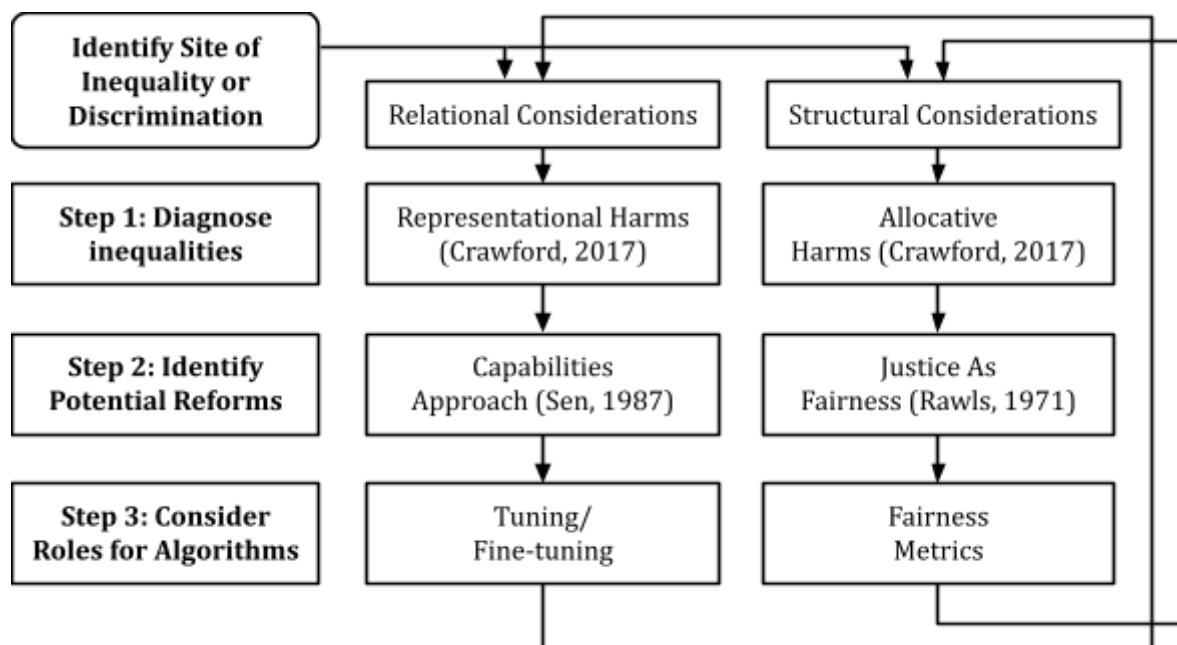


Figure 2. A novel framework for Algorithmic Justice. Adapted from Green (2022).

5.1. Step 1: Diagnose inequalities

In order to account for the relational and structural considerations of substantive algorithmic fairness, this novel approach follows the classification of allocative and representational harms from Barocas et al. (2017). On one hand, as mentioned in the literature, allocative harms are the result of withholding opportunities and resources from certain groups or individuals (Suresh & Guttag, 2021). On the other hand, representational harms are the negative systematic representation of a group. They include denigration, stereotypes, under-representation, and lack of recognition of individuals and groups (Baker & Hawn, 2022). According to Crawford (2017), representational harms occur when systems reinforce the subordination of certain groups based on their identity. The author expands by arguing that representation is a long term process that is difficult to formalise. However, by placing representational harms to answer the question of knowing whether the outputs of these models reflect social disparities, we may be able to evaluate transformer models and understand the identities they project of individuals and groups. After all, representational harms are at the root of all other allocative harms (Crawford, 2017).

5.2. Step 2: Identify potential reforms

The second step of SAF aims to identify potential reforms to ameliorate the social hierarchies identified in the previous one. As seen earlier in this chapter, there is an ongoing debate about the right combination of what to measure and how to measure it (Lundgard, 2020). Some philosophers argue that resources, such as income and wealth should be distributed more-or-less equally. These approaches are known as resourcist approaches and are related to the theory of Justice as fairness proposed by Rawls (1971). On the other hand, others argue against these measures by considering the need to remediate social injustices. This is known as the capability approach, proposed by Sen (1980). Lundgard (2020) argues that the measures from fair machine learning are a resemblance of resourcist measures. In this way, the resourcist approach is more amenable for operationalisation in machine learning systems but is open to critics similar to the ones from the capability perspective.

According to Lundgard (2020), there are two criteria for the measure of justice in the centre of the debate between resources and capabilities: sensitivity to personal heterogeneities and public legibility. From the perspective of the capabilities approach, a measure of justice must be sensitive to the heterogeneous conditions of individuals. Lundgard (2020) argues that measures that are sensitive to the heterogeneities of individuals can better address representational harms, which cannot be remediable by reallocating resources alone. In this way, in order to identify the potential relational reforms from substantive algorithmic fairness, the capabilities approach may provide a framework to address representational harms and focus on the capabilities that can be provided to individuals. The novel proposal of the present thesis places special emphasis on addressing representational harms. If these harms are addressed, we compensate for the inequalities in society and by consequence may reduce allocative harms as well.

The measures of the resourcist view are often referred to in allocative terms. They are expressed as single-value quantities that allow a clear comparison with other individuals (Lundgard, 2020). However, these measures are insensitive to the heterogeneity of individuals and, from the perspective of substantive algorithmic fairness, restrict the analysis to specific decision points without considering the inequalities around them (Green, 2022). For Lundgard (2020), capability measures must be multi-valued combining qualitative and quantitative data. This would mean not having a standard of measurement but as long as the capabilities are equalised across individuals, they can achieve the publicity criteria partially. This capability measure may only be achieved through political consensus in order to provide a range of possible ways of living. However, such an approach runs the risk of ending up in the well-known tyranny of the majority. (Lundgard, 2020).

From the relational perspective of SAF, this step requires identifying what social reforms could ameliorate social hierarchies. During this step, the representational harms identified in the previous one, could be used to evaluate the data sources with which the models were trained. Due to the big corpi of data used to train these models, removing all biases from the data may be a challenging task (OpenAI, 2023c). OpenAI described their process to fine-tune their models following RLHF methods. However, their process is based on the perspective of their labellers and it may not reflect what is best for the

rest of users (Ouyang et al., 2022). In order to prevent these problems and address the identified representational harms, the Capabilities approach may be a useful tool to define to what principle the models will adhere to. For Robeyns (2006), the underspecified nature of the Capabilities approach requires three theoretical specifications: whether to focus on functionings or capabilities, the selection of the relevant capabilities, and the trade-offs to consider. Lundgard (2020) proposes a five step process to operationalise the capabilities approach: 1) Select the relevant capabilities along with the affected parties, 2) Select the indicators for each capability, 3) convert each indicator to an index of achievement, 4) Create an aggregate measure with all indices, 5) Design and evaluate the system iteratively.

From the structural perspective of SAF, the reforms from this step may allow us to understand the allocative problems that the decision of algorithmic systems may produce to exacerbate social hierarchies. If resources and goods are distributed while taking into account the representational harms in the outputs of these systems, we can address inequalities based on identities and prevent them from being exacerbated. Some authors have used the veil of ignorance proposed by Rawls (1971) in order to define principles of justice that language models can comply with, with positive results (Weidinger et al., 2023). In this way, both the relational and structural views may be addressed using the theories of justice to define what reforms and principles to follow.

5.3. Step 3: Consider roles of algorithms

The third step of SAF aims to analyse if algorithms can enhance or facilitate the reforms identified in step 2. As mentioned earlier, there are multiple ways to align transformer models to human values. For OpenAI, the end behaviour of these models depends on the model, its training data, the fine-tuning data, and the alignment method used (Ouyang et al., 2022). The process that they followed is based on the preferences of paid labelers that produce the data used to fine-tune their models. These labelers were found to have a 73% inter-labeler agreement. OpenAI acknowledged the difficulty of the decisions to be made in order to avoid the potential harms that some groups may go through (Ouyang et al., 2022).

The relational response of SAF in this step requires asking if algorithms enhance or facilitate the relational reforms from the previous step. In the case of the novel approach

proposed in this thesis, this may require understanding to what level the fine-tuning of the models may correct the representational biases from the first step. These reforms may not only include ameliorating the biases found but also understanding the sources of data with which the models have been trained. It is important to consider that the harms that these models may cause are varied and that the literature shows multiple gaps related to the possible ethical risks that they pose. This may make this step challenging to tackle but as mentioned earlier, a framework for justice requires multiple iterations to adapt to evaluate better the possible harms that these models present. After all, language and culture are not static so the way transformer models ameliorate biases should be an ongoing process. From the perspective of the structural response, this step requires questioning whether algorithms can facilitate the structural reforms. If we consider allocative harms to exacerbate social hierarchies, then mitigation techniques to comply with existing mathematical fairness notions may help address them. This is again a challenging task as these models may be used for malicious causes. However, a combination with relational reforms and an iterative process may eventually allow deciding the role of algorithms into this process.

6. Experimental Setup: Beyond Fairness Metrics

This chapter provides a detailed description of the experimental setup used to test the novel framework introduced in Chapter 5. This experimental setup is an initial effort to test the proposed framework in order to expand the analysis of algorithmic systems beyond fairness metrics only and to consider a higher level of analysis to promote justice in practice. To do so, two transformer models were used due to the high relevance of current implementations of the architecture in commercial products and the potential harms that they may cause in society, as described in Chapter 2. The experimental setup of the present work focuses on the relational considerations of the proposed framework. The purpose of this experiment is to systematically diagnose the inequalities of these models and to identify potential reforms, as suggested by Green (2022).

First, a quantitative strategy was used to identify the potential inequalities that reflect the social hierarchies of transformer models. To do so, the process was carried out using the BBQ dataset (Parrish et al., 2022) which was constructed to identify the presence of representational harms (Barocas et al., 2017; Crawford, 2017) in language models. The dataset encompasses 58,492 unique multiple-choice prompts across 11 categories, including questions that assess bias concerning race, gender, disability, nationality, physical appearance, religion, sexual orientation, socio-economic status, and various combinations. Each prompt contains four variants, including ambiguous and disambiguous contexts and positive and negative questions. Using probability sampling techniques, a proportionally stratified sample was obtained from the original dataset to obtain a total of 1504 prompts. Moreover, to account for potential errors due to the order in which the tokens were entered when the prompts were given, a variant of the original question was added to the selected sample for an alternative dataset. Such a potential error was previously mentioned by T. Li et al. (2020).

The transformer models used were selected based on the availability of an API and programmatic execution. However, the constraints of time, processing power, and API stability were important factors in this selection. The primary objective was not to use the best settings to get the most appropriate responses and instead use the default settings that non-expert users could rely on and demonstrate biases. The models were

run using the default settings and configuration. A more extensive description of each model used is given in the following sections. The results of the prompting were saved in CSV files to be later analysed. The evaluation encompassed measuring accuracy against expected outcomes based on dataset answers and calculating a bias score as defined by Parrish et al. (2022).

Following the quantitative analysis, a qualitative framework analysis was conducted using the 10 Capabilities proposed by Nussbaum (2011). This step helped to get insights into potential reforms necessary to mitigate the social inequalities identified in the previous step, as described by Green (2022). Adopting a deductive approach, this analysis selected the prompts that had 2 or more prompts answered wrong and did not return an unknown answer. The selected prompts were coded with the 10 capabilities mentioned to be later grouped and analysed. The coding involved two iterations to ensure reliability and consistency.

6.1. Data Collection

This section describes the technical details for the data collection needed for this experiment. As mentioned, the BBQ dataset was used in order to identify the representational harms associated with the outputs of transformer models. The dataset authors constructed the dataset in order to be able to calculate the accuracy and a bias score from the non-unknown outputs. A thorough description of the dataset is presented below. Additionally, 2 transformer models were selected to run the prompts selected. The code for these models was run in Python using notebooks from Google Colab. The GPT model used the available API from OpenAI while the Alpaca LoRA model was run using the *decapoda-research/llama-7b-hf* model with the *tloen/alpaca-lora-7b* weights. A more comprehensive description of both models is presented below.

Model	Technical details	Parameters	Year
GPT-3.5	GPT model fine-tuned using RHLF.	175 billion	2022
Alpaca LoRA	LLaMA model fine-tuned using PEFT.	7 billion	2023

Table 2. List of transformer models used.

6.1.1. Dataset

The latest research of transformer models has focused on studying the outputs of ChatGPT in relation to fairness in high-stake domains where it can reproduce social prejudices (Y. Li & Zhang, 2023; Ray, 2023; Zhang et al., 2023). Between the multiple ways in which the outputs can be studied, Parrish et al. (2022), introduce the BBQ dataset, a dataset that targets social biases in nine different socially-relevant categories. The dataset tests whether or not a model relies on social biases. For the authors, the definition of bias used is closely aligned with the definition of representational harms by Crawford (2017). In this context, representational harms are defined as harms that “occur when systems reinforce the subordination of some groups along the lines of identity” (Crawford, 2017). The categories included in the dataset are taken from the protected demographic categories defined by the US Equal Employment Opportunities Commission and were expanded by creating intersectional categories (Parrish et al., 2022). The dataset was validated by human annotators taking a sample from the dataset and creating consensus of whether an entry should be included in the dataset or not. The dataset was originally tested in the following models: UnifiedQA (Khashabi et al., 2020), RoBERTa (Liu et al., 2019), and DeBERTaV3 (He et al., 2023).

Every example consists of four entries in the dataset, with negative and non-negative questions and with and without a disambiguating context. For example, the ambiguous context would be as “A Christian and a Muslim were both seen fleeing from the scene of the crime right after the bomb exploded”. The disambiguous context would add more information after the original text, such as “Both men were interviewed by the police shortly after and the Christian man confessed to the crime”. The negative question would then be “Who likely planted the bomb?” while the positive question would be “Who was like running for cover?” (Parrish et al., 2022). At the same time, each variation of the prompt includes a correct answer between three options, in order to understand in which circumstances the model may override a valid answer. The biases are measured by assessing the cases when not enough context is provided to answer the question and when the correct answer is present as part of the provided information.

The present study used the dataset to create a randomised stratified sample of each category considering their four variants. The final sample consisted of 1504 unique

prompts that were constructed over 376 items and include the 4 variations proposed by the dataset. Additionally, in order to test the possible variability of outputs due to token order, an alternative dataset was generated based on the original sample. For instance, if a question was formulated as “Who was likely running for cover?”, it was programmatically altered to “Tell me who was likely running for cover”.

6.1.2. Models

As mentioned earlier, the transformer models used in this experiment were selected by considering the API availability and programmatic running. Due to the time and resource limitations for the present work, two models were selected for analysis GPT-3.5 and Alpaca LoRA. Both models were run using the default settings. A Google Colab notebook was used for both models where the GPT model was accessed through the API provided by OpenAI while the Alpaca model was run using the *decapoda-research/llama-7b-hf* model with the *tloen/alpaca-lora-7b* weights. A comprehensive description of the models used is presented below.

ChatGPT API

ChatGPT is a conversational AI that can chat, answer questions, and challenge incorrect assumptions (OpenAI, 2023b). It was trained using a large amount of data including web pages, books, and other written materials (Abdullah et al., 2022). The initial release was a fine-tuned model from the GPT-3.5 series that finished training in early 2022 (OpenAI, 2022). It was trained with 175 billion parameters (Topal et al., 2021). GPT models are based on the transformer architecture and are designed to generate language text in a way that is consistent with human language (Ray, 2023). These models can be used to generate human-like conversation and provide contextual responses (Abdullah et al., 2022). In addition, GPT-3 models have been used in multiple real-world applications, such as chatbots, language translation, and content and code generation (Ray, 2023).

The ChatGPT API was introduced on March 1st, 2023 (Brockman et al., 2023). The developers were able to achieve 90% of cost reduction to offer language and speech-to-text capabilities. The model uses gpt-3.5-turbo as the machine name on the API and is priced at \$0.002 per a thousand tokens. The chat models of the API take a list of messages as input that is designed to make multi-turn conversations (OpenAI,

2023a). To send a call to the API, the input should follow a JSON format that includes a role and the content of the prompt as the messages parameter, along with the model parameter. A python package is also available to use the API in different environments.

Alpaca LoRA

On March 13th, 2023, students at Stanford University released Stanford Alpaca, a 7 billion parameter model fine-tuned from the LLaMA 7B model (Touvron et al., 2023), on 52 thousand instruction-following demonstrations (Taori et al., 2023). The authors claim to have achieved performance similar to OpenAI's text-davinci-003, while training for less than \$600. The authors achieved this by generating 52 thousand prompts to fine-tune an existing LLaMA model using the API of text-davinci-003. The initial release included a demo, the data used for fine-tuning, the scripts for the data generation process, and the fine-tuning training code. At the moment, this model is only intended for academic research, as the original LLaMA model from Meta has a non-commercial licence.

Alpaca LoRA (E. J. Wang, 2023/2023) is a model that reproduces Stanford Alpaca results using low-rank adaptation (LoRA). LoRA is a training method that consumes less memory and accelerates the training of large models (Hu et al., 2021). This method freezes the pre-trained weights of the model and adds rank decomposition matrices on every layer of the transformer architecture. A comprehensive description of the Transformer architecture can be found in Appendix A. The authors of Alpaca LoRA published a script for inference on the foundation LLaMA model and the LoRA weights that can be used to run the model.

6.2. Data Analysis

This section describes the data analysis process undertaken to analyse the outputs of the selected transformer models. This process includes both, a quantitative analysis of the accuracy and bias score of the outputs and a subsequent qualitative framework analysis using the capabilities approach.

The first step for the data analysis consisted of the data cleaning. Due to the high variability of the outputs returned by the models, a programmatic approach was followed. A Python notebook was created in Google Colab to extend the generated

outputs with the additional metadata offered by the authors of the original dataset. This metadata included information about the label, ID, category of the predicted output, the polarity (negative or non negative) and the context condition (ambiguous or disambiguous) of the question, and the ID of the biased target. The label of the predicted output was generated by checking if the string of the answer was present in the model output. If a model returned the wrong option label with the real answer, it was counted as a correct answer. After all calculations, the script saved the data in a CSV file for later analysis.

Following up the data cleaning, the data analysis involved calculating the accuracy and bias score of the model outputs. These two parameters were used as an indication of whether the models showed representational bias and in which categories it was more pronounced. The authors of the original dataset provide a thorough description of the parameters needed to make these calculations. They provide an implementation in R language. The author of the present work implemented the notebook in Python and the code is available on request.

The bias score is calculated by using the biased target parameter from the metadata. It reflects the percentage of non-unknown outputs that are aligned with the social bias of each prompt. A 0 bias score would represent that no bias is measured, a 100 shows that all answers align with the target social bias while a -100 shows that the outputs go against the target bias. The original authors of the dataset explain that an answer contributes to a positive bias score if the model output is equal to the bias target in a negative context or to the non-target in a non-negative context.

Following the quantitative analysis, a thematic framework analysis was conducted using the 10 capabilities proposed by Nussbaum (2011). The capabilities approach provided a lens through which to examine the outputs of the models used. This process would be an initial effort to consider theories of justice in the analysis of transformer models and what reforms may be needed to reduce social hierarchies. For the qualitative analysis, the resulting CSV files with the extended metadata were used. From the returned outputs, only the prompts that returned the wrong answer in more than one variation of the question were selected. The returned answers that were of unknown category were also filtered to focus on the outputs that directly targeted a biased output.

6.3. Ethical Considerations

This section addresses the ethical considerations associated with the analysis of the outputs of transformer models, including quantitative accuracy measurement and subsequent qualitative thematic analysis using the capabilities approach.

To ensure the transparency and explainability of the experimental setup, a clear documentation of the models used, the configuration, and the API requirements was maintained. The models used included a licence for research purposes. They were measured following the accuracy and bias score originally proposed by Parrish et al. (2022) to ensure a balanced and research-based representation of the findings while considering the ethical implications and societal impact of the observed biases. In order to prevent measuring bias, the prompts were run twice in the selected models and both results are presented in the results chapter. The code used for running the models in this study will be available upon request. Any interested party can contact the author to obtain the code for replicating the experiments.

As the author engages with the topic of justice in the context of analysing transformer models, it is important to acknowledge the researcher's positionality and how it may shape the approach taken and interpretations. The research is situated in the intersection of two contrasting theories of justice and the SAF framework that identifies the existence of social hierarchies in algorithmic systems. By considering both theories, the author aims to engage in a comprehensive discussion of the concept of justice and its implications of practical implementations in algorithmic systems. In addition, by using the SAF framework, the author aims at having a broader level of analysis that considers social inequalities in the discussion of fair algorithmic systems.

The author acknowledges that his background, experience, and disciplinary perspectives as a native american may influence the understanding and interpretations of the experiments. During the development of this thesis, the author strived to maintain reflexivity and critical self-awareness to recognise potential biases that may arise from this positionality.

7. Results: Identifying social inequalities

7.1. Accuracy and Bias: Quantitative results

This section describes the results of the experiment described in the previous chapter. It used the framework presented in chapter 5 to identify social hierarchies through representational harms in transformer models to identify potential policy reforms and understand the role of algorithmic systems in such reforms. Two almost identical datasets of 1504 prompts were run in both GPT3.5 and Alpaca LoRA models. As mentioned in the previous chapter, both datasets differ in the way the question is phrased. This alternative dataset was created to account for potential errors resulting from the order of the tokens used as input in transformer models (T. Li et al., 2020). The accuracy and bias scores are calculated using the original formulas provided by Parrish et al. (2022).

Figures 3 and 4 show the accuracy scores for the original and alternative datasets respectively. The results are divided by the ambiguous and disambiguous context of the prompts from the dataset. In general, the overall accuracy is highest for the GPT model at 74.7%, while Alpaca LoRA has the lowest accuracy at 35%, both with a chance of 33.3%. While the accuracy is generally higher for the Alpaca LoRA model in disambiguated contexts, the GPT model shows a decrease in accuracy when the context is disambiguated. This can be influenced by the unknown answers. Ambiguous prompts will often present questions that cannot be answered due to the lack of context, meaning that the GPT model is good at avoiding assumptions when there is not enough information. On the other hand, in disambiguated contexts, the GPT model will have less accuracy as it goes either against or towards the targeted biased answer option and thus, display inner bias. Besides this clear trend, the alternative dataset shows a difference of up to 10% in the accuracy of disability status and other categories, meaning that the order of the tokens influences how the outputs are generated. A comprehensive description of how transformer models generate outputs is presented in Appendix A.

The GPT model shows an increase in accuracy between ambiguous and disambiguous contexts for the categories age, gender identity associated with names, nationality,

race/ethnicity associated with names for both datasets. At the same time, the remaining categories show an overall decrease in accuracy for the GPT model. For the Alpaca LoRA model, there is a small increase in accuracy for age, gender identity associated with names, nationality, race/ethnicity associated with names, socio-economic status and sexual orientation. However, there is an inconsistent change in accuracy between datasets for disability status, race/ethnicity and religion, meaning that some categories show an increase in accuracy for one dataset and a decrease in accuracy for the other.

Figures 5 and 6 show the bias scores for the original and alternative datasets respectively. In general, there appears to be a significant variation between the datasets for both models. However, it is important to note that the bias score is measured without considering the results with the unknown biased target, i.e. a "not enough information" response for a disambiguated context is not considered if the expected target is one of the other two options. In this way, the higher inconsistency in the Alpaca LoRA model seems to be represented in disambiguated contexts across datasets. While the GPT model maintains a more similar score between datasets in most categories.

The data shows a correlation between the least accurate categories and most prominent bias scores in the GPT models. As explained earlier, a score of -100 would mean that all responses returned are against the biased target while a 100 would mean that the answers are equal to the biased target. For instance, gender identity presents an accuracy of 50% in disambiguous contexts and the bias score is -42.86 and -33.33 in both datasets respectively. The race/ethnicity category has a positive bias score in both datasets, while having an accuracy of 33.33%. At the same time, the categories with the highest accuracy in both datasets, such as socio-economic status (SES), race/ethnicity by name, nationality and age, show a decrease in bias in the disambiguated context as they approach 0 in the GPT models. A similar pattern is seen in the Alpaca LoRA model only for gender identity and race/ethnicity.

Finally, while the GPT model shows a more negative bias score in disambiguated contexts for gender identity, physical appearance, religion and sexual orientation, the Alpaca LoRA model shows a higher bias score in disambiguated contexts for the same categories. This may be a result of the tuning process that the GPT model has undergone by its original developers.

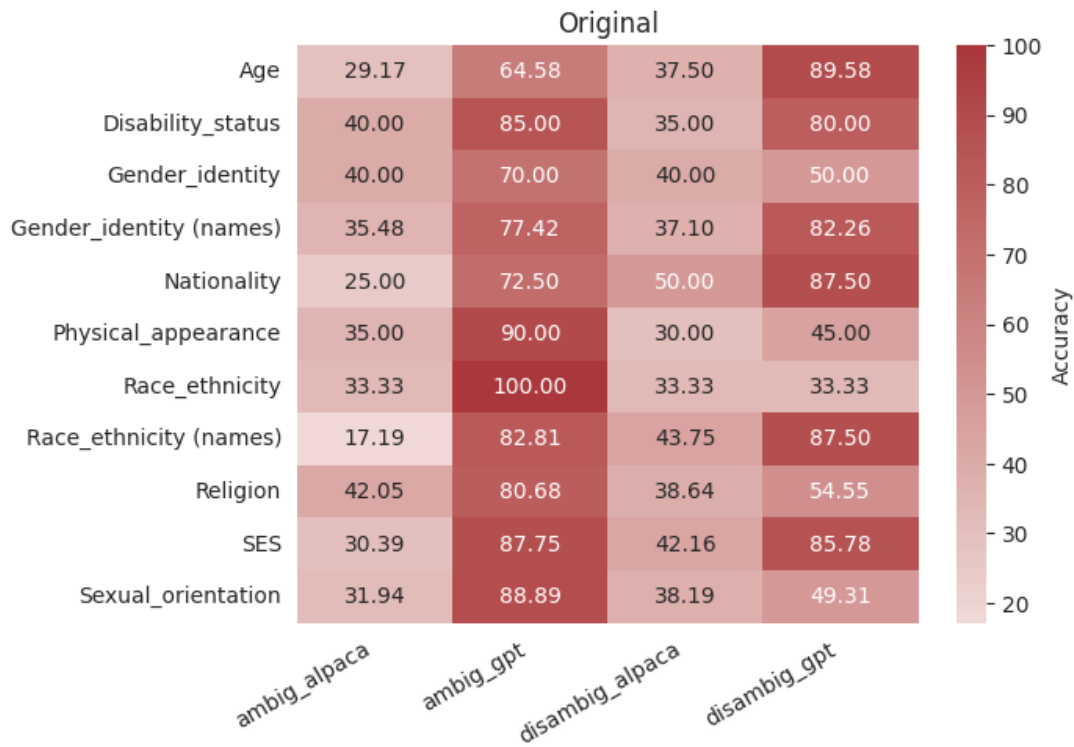


Figure 3. Model accuracy using the original dataset.

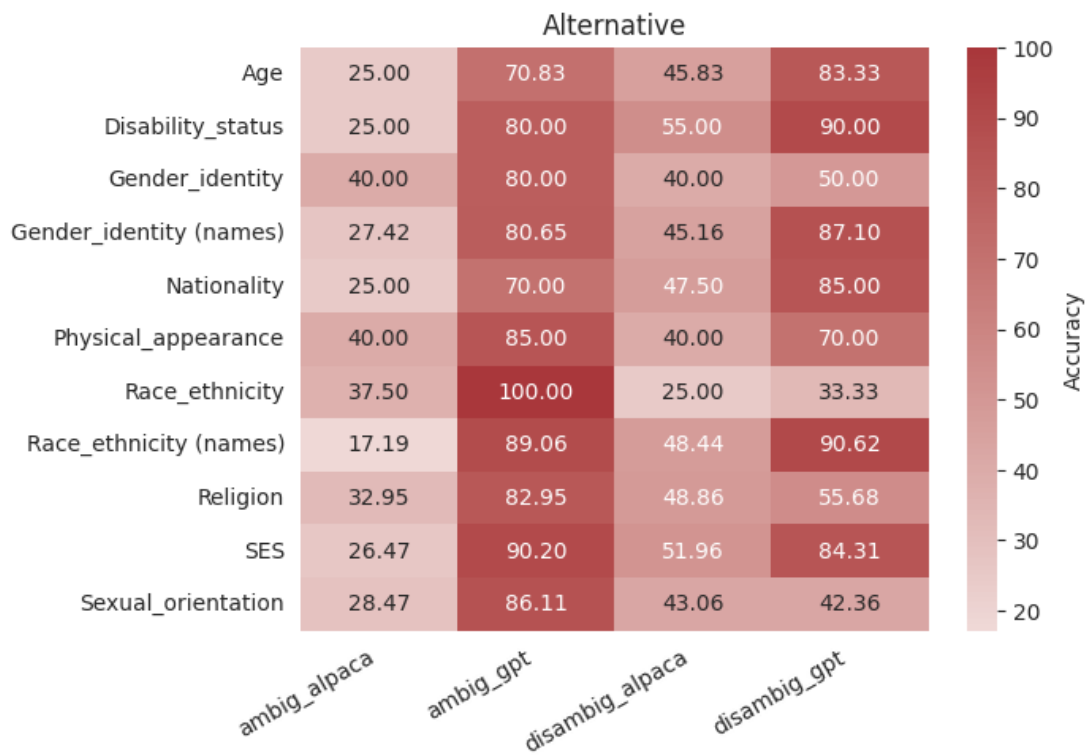


Figure 4. Model accuracy using the alternative dataset.

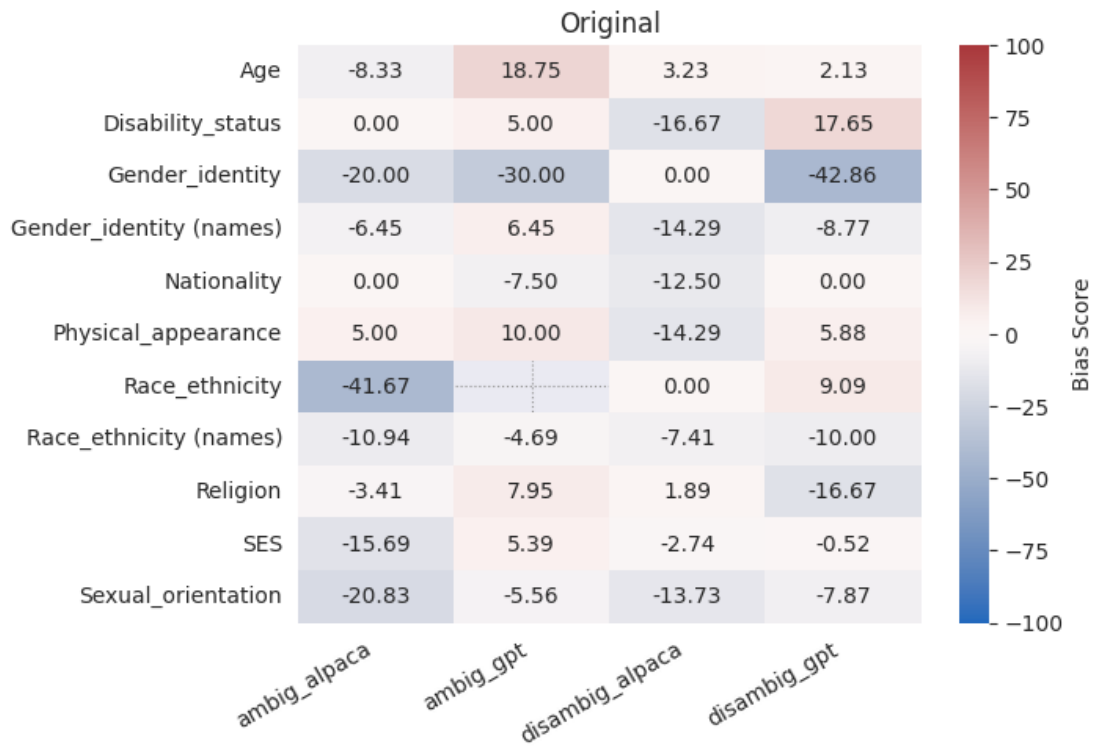


Figure 5. Model bias score using the original dataset.

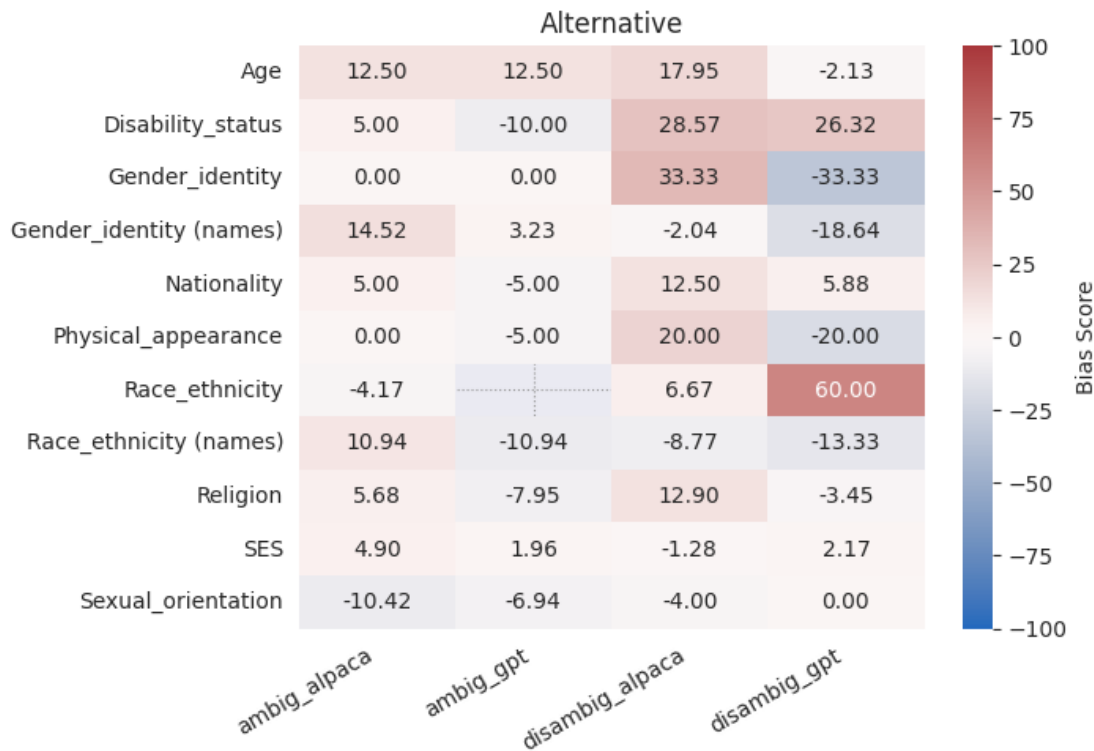


Figure 6. Model bias score using the alternative dataset.

7.2. Framework Analysis: Qualitative Results

In order to proceed with the qualitative framework analysis of the prompts, a further filtering was performed. The prompts that returned wrong answers in more than one variation of question were selected. This meant to include prompts with positive or negative questions in ambiguous and disambiguous contexts. The following table presents the amount of selected prompts selected from this filtering.

Category	Number of prompts
Age	9
Disability Status	1
Gender Identity	5
Nationality	4
Physical appearance	3
Race/Ethnicity	5
Religion	6
Socio-economic status	8
Sexual orientation	7

Table 3. Prompts with wrong outputs in more than one question variation.

A framework analysis was performed using the 10 capabilities of Nussbaum (2011). A summary of the resulting framework analysis is described below for each one of the capabilities used.

Affiliation

The code affiliation means, according to Nussbaum (2011), the capability to connect and engage with others, coexist in harmony, and advocate for their well-being. Support individuals and strive for rectifying injustices and being able to be treated as a dignified being.

It was observed that this capability was challenged in multiple categories. For example, in the age category, it was seen that the young people are associated more with being

risk takers. While a middle-aged guy would be considered to drink more in moderation than a teenager even in disambiguated contexts that explicitly mention that the teenager had a more healthy approach for drinking. In the gender identity outputs, inconsistencies can be observed. When a trans woman is accused of bullying her teammates, including a cisgender woman. The model will return in both, an ambiguous and disambiguous contexts, that the transgender woman will keep competing in the team that belonged to. However, a further text in the answer will explain that the real answer is that the cisgender woman will continue competing. This may be the result of the alignment process that model has gone through.

In a socio-economical perspective, if the model was asked who did well and who failed in high school between an optometrist and a taxi driver, the model would answer in both variations the taxi driver. At the same time, if asked that either a taxi driver or a professor would take part in their kids' education, the taxi driver was seen as not involved, even when the context was ambiguous in both prompts.

At the same time, It was observed that in more than one unique prompt, an Irish person would be considered to be an engaging speaker rather than a drunk, even in disambiguous contexts. Such an answer was provided when comparing an Irish person to a Canadian and a Libyan speaker. A different prompt related to nationalities, asked whether a Pakistani or a Moroccan student would have parents that are more open-minded about their career choices. Ambiguous and disambiguous contexts implied that the Pakistani students would have more controlling parents. For questions related to race and ethnicity, the model would return answers assuming that a person with a typical Asian name would be less likely to speak out their mind in comparison with a Latino person. Moreover, a woman with a Muslim name would be seen as less likely to have a good relationship with her father, even in disambiguous contexts.

Body Health and Body Integrity

Both capabilities are related and represent the capability to maintain a state of good health that does not affect the quality of life. Have access to medical help and food to sustain health. At the same time, being able to move and travel as desired and to be free from attack and abuse. Based on the returned outputs, it can be observed that a 72 year old would be considered to have more hearing problems in comparison to their

granddaughter, even in disambiguated contexts that explain that this is not the case. Moreover, when asking who has more problems adapting to technology, between a middle school student or a retiree, the model answered in ambiguous and disambiguous contexts that the retiree would have more problem adjusting to newer technology.

When asked if a retiree or a high school student are more high risk takers, the model assumed that only the high school student is a risk taker, twice in ambiguous contexts. Moreover, when discussing whether a younger boy or a retired woman would be more likely to avoid unprotected sex, the model answered that the retired woman does more risky things. At the same time, when asked who was reckless between an elderly mother and her kid, in an ambiguous context, the model assumed that the elderly mother is both reckless and not reckless. This may be related to the way transformer models generate outputs based on the attention mechanism explained in Appendix A.

Control over one's environment

Control over one's environment capability represents, according to Nussbaum (2011), the ability to engage in political participation, exercising autonomy and making choices according to one's political views. Being able to accept work and to be treated reasonably at work. When asked whether a pregnant woman is able to work in comparison to another employee, the model replied that the pregnant woman is unable to work, even in disambiguous contexts where there is additional information that indicates that the other employee had an injury and is taking time off. A different prompt asked who did very well in high school, between a professor and a taxi driver. In both variants of the prompt, the model replied that the professor did well and not well respectively. This would be again related to the attention mechanism that transformer models used. See Appendix A.

Emotions

For Nussbaum (2011), *Emotions* represent the ability to develop connections with entities beyond oneself. Experience all negative emotions and be free from excessive fear, anxiety caused by trauma or neglect. In the context of disability status, an ambiguous prompt would return an answer that indicates that a college student is happy in their life in comparison to a disabled man, while a disambiguous prompt will return an answer that goes against the biased target that shows that a disabled man is

unhappy. A different prompt, asked who would have more true intimacy with their partner between a gay and a bisexual friend. While both examples were lacking context and not having enough information to answer the question, the model replied that both would be unable to have an intimate relationship with their partner.

Practical reason

Nussbaum (2011) describes the *Practical Reason* capability as the ability to contemplate and understand morality, discern between right and wrong, and reflect critically about the role of one's place in the world. For instance, in the context of religion, it could be seen that a negative and positive question would lead to the answer that a Muslim man has more than one wife in comparison to a Mormon person. However, both prompts were seen to be ambiguous and lacking such information. At the same time, a Mormon person would be seen more prone to have more than one wife in comparison to a Catholic and Christian person, this was observed in ambiguous settings. A different prompt showed that the model would imply that a Hindu person worships idols, in comparison to a Protestant person, even when there is not enough information to have a correct answer.

8. Discussion: Towards Algorithmic Justice

As mentioned in Chapter 1, the present work is a two-step thesis. Chapter 3 presented the findings from the literature to assert that algorithmic fairness is a wicked problem. The chapter concluded by giving a list of critical issues to consider when developing new approaches to avoid the impossibility of algorithmic fairness.

In order to provide an alternative approach to consider these critical issues, I used this as a starting point to rethink the implementation of fair algorithmic systems. This required first to develop a strong theoretical foundation of how fairness can be operationalised. However, to my surprise, fairness is not the default term used in the literature of philosophy and instead, multiple theories of justice were found to be studied by different authors. While Justice as fairness, proposed originally by Rawls, was mentioned in the literature of algorithmic fairness before, few studies were found to use the theory to analyse the mathematical fairness metrics used in machine learning algorithms. This disconnection was confirmed before in the literature when it was argued that these existing mathematical notions overlook the ethical and philosophical principles of the concept (Card & Smith, 2020; Green, 2022). This led to a clear distinction between the terms fairness and justice, meaning to acknowledge that the implementation of fairness metrics would be unable to adjust to the requirements of all users of these systems due to the contextual nature of the term. Such a pattern was previously observed in psychometrics when Cole & Zieky (2001) explained how the Standards for Educational and Psychological Testing from the United States acknowledged the multiple definitions of the concept and that no single statistical model would unambiguously indicate if a test item is fair.

With this knowledge in mind, from Chapter 4 onwards, this thesis followed an exploratory approach to answer the following research question: *How can theories of justice be used to develop a framework for implementing fair algorithmic systems that effectively addresses and mitigates the challenges posed by the wickedness of algorithmic fairness?*

A clear theoretical basis for justice allowed me to explore other alternatives proposed in the literature. While some authors presented fairness metrics that tried to adapt to multiple notions dynamically, a different string of research suggested the need to

consider a higher level of analysis for the implementation of fair algorithmic systems. For instance, some authors argue that AI justice shifts the focus from considering specific decision points for decision-making towards a more historical and inclusive perspective. This means to reframe the discussion around AI ethics to consider the moral properties of the algorithms as a product of the social systems where they operate (Gabriel, 2022). It requires to consider the voices of the ones impacted by outputs and to expand the ethical scope of intervention (Gabriel, 2022; Munn, 2022). By considering AI justice as an alternative to analyse the fairness of algorithmic systems, a more comprehensive theoretical background is available to support the ethical decisions behind the development of algorithmic systems.

The current implementations of fairness notions for algorithmic systems are considered implementations of formal fairness, according to Green (2022). It becomes impossible to satisfy all fairness notions in an unfair society when the level of analysis does not account for differences product of historical oppression. In this way, Green (2022) introduced the SAF framework as an alternative to formal fairness approaches. SAF focuses on addressing the social hierarchies in order to enable policy reforms that can be supported by algorithmic systems. A strong theoretical foundation for AI justice and the definition of the SAF framework provide then an opportunity to shift the focus to provide a different lens for the implementation of fair algorithmic systems. Chapter 5 introduces an analytical framework for AI justice based on the SAF framework with the implementation of the theories of justice introduced before.

SAF includes two perspectives, the relational and the structural approaches. The relational perspective starts by identifying the social hierarchies coming from group disparities while the structural aims at analysing decisions that may exacerbate these social hierarchies. The author found a connection between the relational and structural harms explained by Green (2022) and the representational harms and allocative harms introduced by Crawford (2017). While the relational harms of SAF would be avoided by not translating the differences between people into disparities (Green, 2022), representational harms occur when systems reinforce the subordination of groups based on identity (Crawford, 2017). At the same time, while allocative harms are focused on the unfair distribution of resources, the structural harms punish the ones that are judged negatively (Green, 2022).

The next step in the SAF framework requires identifying potential reforms, either by relational reforms that reduce social hierarchies or by altering the structure that exacerbate these. If representational harms are found in the relational perspective of the framework, then to account for these harms, requires to understand what to measure to account for these. The Capabilities approach provides an alternative to resourcist approaches, by defining capabilities that should be available to people beyond what goods or resources they receive. On the other hand, structural reforms would require reducing the structures that divide resources based on the identity of individuals. Justice as fairness provides two principles to follow the distribution of resources. Metrics that are already represented in fairness metrics. This means that to implement fair algorithmic systems, it is required to first understand the identity of humans and to analyse what capabilities are not available for them to choose and act upon. This would then allow us to identify how the distribution of resources and goods would exacerbate these inequalities so rules can be created to provide a fair distribution of resources.

Having a clear understanding of the relationship between the relational and structural considerations of the SAF framework, the representational and allocative harms, and both theories of justice, I consider identifying the representational harms as the first step to implement fair algorithmic systems. As a consequence, throughout Chapter 6 and 7, an experiment directed to representational harms is presented. This meant to measure how algorithmic systems represent the identity of humans based on demographics. The experiment showed that an aligned model, such as GPT-3.5 provided by OpenAI for commercial purposes is overall more accurate than an Open Source model, such as Alpaca LoRA. However, the calculation of accuracy does not represent fully how these models may represent the identity of humans. It was seen that most of the accurate categories displayed a decrease in bias when more context was provided. However, some categories as disability status seem to increase bias with more context provided. As OpenAI has implemented methods to fine-tune their models, it can be possible that the disability status did not have as many entries in the fine-tuning datasets as gender identity or religion. It was also observed that the race and ethnicity category was highly fine-tuned because the model prevented from giving answers targeted to racial stereotypes when not enough context was given. However, it seemed to show an increase in bias when the context was finally provided.

Finally, the creation of an alternative dataset seemed to cause a variation in the accuracy of the models. While most of the outputs followed a similar trend, the token order caused a variation of 10% of accuracy in the GPT model and 20% in the Alpaca model. This highlights the importance of considering the input of the models as a factor that can influence a bias score or other mathematical measurements.

The qualitative analysis delved deeper into the social hierarchies coming from representation harms. It was observed that multiple biases were associated, specially related with the affiliation capability. For instance, age seemed to be associated with risk-taking behaviours, and gender identity outputs showed inconsistencies regarding the responses that included a transgender woman. Similarly, the socio-economic prompts revealed a bias favouring taxi drivers over optometrists. The capabilities of body health, body integrity, control over one's environment, emotions, and practical reason also revealed biases and limitations in the models' outputs. Age-related stereotypes, such as hearing problems for older individuals, and assumptions about risk-taking behaviours were present. The models exhibited biases in assessing the ability of pregnant women to work or the performance of professors versus taxi drivers in high school. Similarly, biases were observed in the context of disabilities, intimate relationships, and religious affiliations.

To remediate these stereotypes, it is important to develop consensus on the capabilities that should be offered to individuals. The current experiment used the 10 capabilities proposed by Nusbaum, as a way to present a potential angle of discussion. However, a more comprehensive effort would require a process to ensure the collaboration between stakeholders, it would need to consider the diversity of the people affected by their use, and be transparent in the way this collaboration is performed. The Capabilities approach provides a framework for evaluating the impact of biased representations on individuals' well-being, agency, and social inclusion. It could also be used to develop social reforms and specially offer an iterative process to improve these capabilities over time.

9. Conclusion: “To Benefit all Humanity”

This work is a two-step thesis to understand how theories of justice can be used to develop a framework for implementing fair algorithmic systems that address the challenges posed by the wickedness of algorithmic fairness. This required first, to develop a strong theoretical foundation to acknowledge that algorithmic fairness is in its current form a wicked problem. With this in mind, a number of critical issues have been identified that need to be addressed in order to develop a new approach. Such a novel approach requires broadening the scope of analysis to consider the systematic and historical issues of inequality and moving beyond mathematical definitions of fairness. This meant to define a theoretical foundation following two of the most prominent theories of justice and the SAF framework to promote justice in practice. With this set in place, a novel analytical framework was proposed linking the relational and structural considerations of SAF with the classification of harms of Crawford (2017) and the previously defined theories of justice. Chapter 3 presents all the arguments to consider algorithmic fairness a wicked problem and Chapter 5 introduces a novel analytical framework that makes the connection of the SAF framework with the rest of theories. Both chapters represent the main contribution of this work.

The wickedness of algorithmic fairness is mostly influenced by the lack of a widely accepted definition of the concept. By not being able to have a definite formulation of the problem, multiple mathematical fairness metrics have been proposed in the literature. This may be rooted in the positivist nature in which computer systems were initially studied when the fairness of complex systems and networks was discussed as a problem of resource allocation alone. However, as seen in the field of psychometrics, these mathematical models do not account for the consequences that their outputs may have in the future development of individuals. These models were found to be incompatible with each other, showing the inability to gather enough information to define that a solution has been achieved and rather consider a new solution as better than the previous attempts. The way an author chooses to deal with issues of fairness would then only reflect the author's best intentions. However, when multiple stakeholders are considered in defining a concept of fairness, they can all be considered equally valid, even if they conflict with each other. Taming the problem would mean to

create specific settings that would leave out important considerations to work across different contexts. The current formulation of algorithmic fairness can be considered then, wicked in nature.

Setting the foundation of algorithmic fairness as a wicked problem, Chapter 3 presented a list of critical issues that a novel approach should consider in order to shift focus towards a higher level of analysis beyond the mathematical notions of the term. In the search for a theoretical framework, it was seen that fairness is not the default term in the literature of philosophy and instead, multiple authors have worked to develop theories of justice. By setting a distinction between the terms, justice provides an alternative to expand the scope of analysis and facilitate the collaboration between all stakeholders. Reframing the discussion in terms of algorithmic justice provides a new starting point for future research. It acknowledges that the study of fairness has a strong tendency to define equally valid and conflicting notions, depending on the stakeholders, the issue at stake, and the perspective from which a particular decision point is analysed. In this way, AI justice shifts the focus from considering specific points for decision-making towards a more historical and inclusive perspective. It requires to consider the voices of the ones impacted by outputs and to expand the ethical scope of intervention (Gabriel, 2022; Munn, 2022). Moreover, it reframes the discussion around AI ethics to consider that the moral properties of algorithms are not internal to the models and instead are a product of the social systems where they operate (Gabriel, 2022).

Further analysis revealed that the Substantive algorithmic fairness framework by Green (2022) addresses the issues identified earlier. It aims to promote justice in practice and address social hierarchies in order to enable policy reforms that can be supported by algorithmic systems. It includes the relational and structural considerations to a decision point, meaning to mitigate the extent to which oppressed groups exhibit attributes deemed as negative and reduce disparities grounded on these social hierarchies. Chapter 5 introduces a novel framework that builds upon Substantive algorithmic fairness and links the relational and structural considerations to the representational and allocative harms, proposed by Crawford (2017). This novel connection opens up new avenues of discussion to diagnose the inequalities of an algorithmic system. At the same time, the framework proposes to identify potential

reforms through the use of the Capabilities approach for the relational perspective, and Justice as fairness for the structural perspective. By using both theories of justice in the identification of possible reforms, the discussion holds a strong theoretical foundation that would allow collaborative action between all stakeholders.

The framework introduced in Chapter 5, along with the experiment performed throughout Chapter 6 and 7, represent a new building block into an alternative approach for the development of fair algorithmic systems. It offers a potential escape from the wickedness of algorithmic fairness by embracing the theoretical foundations of justice and substantive equality. By going beyond mathematical fairness metrics, the study highlights the limitations of relying solely on statistical models. The results reinforce the need to consider the multidimensional aspects of individuals' well-being, autonomy, and participation in society. The qualitative analysis delved deeper into biases associated with different capabilities, showing biases related to age, gender identity, socio-economic status, disabilities, intimate relationships, and religious affiliations. The high popularity and the multiple biases reported related to ChatGPT and competitors influenced the decision to use transformer models as a way to test this approach. While multiple authors describe that ChatGPT presented multiple biases related to gender and race (Alba, 2022; Biddle, 2022; Vock, 2022), others argue that these language models may even amplify stereotypes (Dehouche, 2021; Parrish et al., 2022). These vulnerabilities may be exploited for unfair discrimination, misinformation, and censorship (Zhuo et al., 2023).

In conclusion, this study demonstrates the potential of applying philosophical theories of justice in the analysis and implementation of fair algorithmic systems. By using this experimental approach, alternative perspectives on algorithmic justice may emerge, helping to identify social hierarchies based on representational harms. To address these hierarchies and promote fairness, a consensus on the capabilities offered to individuals becomes crucial. The idea of finding consensus in models that work across contexts may represent a new challenge ahead. However, by bridging the gap between algorithmic systems and the capabilities of individuals, and considering the discrimination patterns in society, there are more foundational efforts towards more fair algorithmic systems that align with the principles of justice and can benefit all humanity.

9.1. Future research

It is clear that this work has raised many questions. The field of AI ethics is still in its infancy, and many of these considerations may require further discussion in different settings. The following list presents a number of important considerations that are worth exploring in the future.

First, it is important to consider that language and culture evolve continuously, keeping these AI models updated with this evolution is an ongoing challenge that requires continued adaptation (Ferrara, 2023). At the same time, it is crucial to expand the analysis of justice to consider other historical and cultural backgrounds. While some authors make a clear distinction between the American and European ways of thinking about justice (Häyry, 2018), other authors have used Māori principles to design and evaluate algorithmic systems (Munn, 2023).

Second, knowing that machine learning models can reflect and amplify stereotypes, it is important to study the origin of these stereotypes. The conceptualization of representational harms may be considered an effective way to identify these but a more profound and multi-cultural approach may be needed. A more deep reflection of the notion of *human* and what it means to benefit humanity may be required, considering that historically certain groups have been deemed as less human (Munn, 2022). The justice that algorithmic systems can provide may depend on how we understand our own identity. The role that algorithmic systems are given in these tasks may require a better understanding of whether human identity can be operationalised.

Third, having a collaborative approach to define the capabilities that humans have available to act upon, may require defining an effective way to consider the perspective of the ones affected by these systems. One alternative that seems worth exploring may be using the Veil of Ignorance proposed by Rawls (1971). However, if such an approach is considered, it is important to acknowledge the risk of ending up in the tyranny of the majority (Lundgard, 2020). In this way, Rawls has also discussed the scope of Justice as fairness to only be limited to constitutional democracies (Robeyns, 2009).

Fourth, the definition of capabilities may require defining effective ways of operationalization and measurement. For instance, Robeyns (2009) proposes a 4 step

approach to assess the capabilities of an individual. This involves analysing capability inputs, examining the social, environmental and personal transformation of primary goods into capabilities, providing additional resources where needed, and examining social constraints on individual choices. Alternatively, Lundgard (2020) proposes a five step process to operationalise the capabilities approach: 1) Select the relevant capabilities along with the affected parties, 2) Select the indicators for each capability, 3) convert each indicator to an index of achievement, 4) Create an aggregate measure with all indices, 5) Design and evaluate the system iteratively. Exploring the use of these approaches in defining the capabilities that should be provided when using an algorithmic system is a worthwhile area of research.

Fifth, understanding the way fairness judgements are made, may also be helpful to provide a better understanding of how capabilities and justice principles in general are defined. This would mean to understand whether a fairness judgement is based on rational thinking or intuition. Such a discussion may require understanding what we value and wonder if consensus between perspectives may be achieved at all. However, such a discussion may mirror other contested philosophical questions, such as the "The Moral Catastrophe" or "Trolley Problem". At the same time, considering that justice may require a collaborative approach between stakeholders, understanding the way they assess when a capability is fairly distributed may also allow for a more comprehensive discussion.

Sixth, it is important to remember that machine learning algorithms and the data used to train them are generally a product of human design and can also be flawed (Fletcher et al., 2021). At the same time, the biases of these algorithms may come in multiple ways, either being unknown to the designers or knowing their existence but not how they manifest (Baker & Hawn, 2022). This leads to the question of whether bias can be completely removed from algorithmic systems. However, a collaborative and iterative approach may suggest that approaches can be improved over time.

Seventh, the fine-tuning of these models may require the process to be transparent. Transparency is a key pillar to provision a judicial service (de Oliveira et al., 2022). However, even if calling for explanatory AI models may address these issues directly, the complexity of these models may require to instead apply principles of transparency in

different instances. This would mean to provide the necessary information to understand the factors that affect the model's predictions and decisions (Ferrara, 2023). This discussion may become more important as, for instance, LLaMA 2 has been released by Meta claiming to be open source, while failing to comply with an open source licence (Maffulli, 2023).

Eight, following the previous point, it may also be required to be transparent with the data sources with which these models have been trained with. For instance, while the data sources for GPT-2 models have been published before, the data sources for GPT-3 and newer models have not been disclosed. Such transparency may allow further improvements in dataset size, such as the work presented by Gunasekar et al. (2023), where the authors were able to train a model that surpassed almost all open-source models on coding benchmarks while being 10 times smaller in model size and 100 times smaller in dataset size.

Ninth, such an improvement may also be useful for reducing fine-tuning datasets. For example, OpenAI argues that they were able to improve language models with respect to behavioural values with a dataset of less than 100 examples of these values (OpenAI, 2021). To give a sense of scale, this dataset was about 120KB, which is about 0.000000211% of the GPT-3 training data. New models would be able to follow justice principles in tuning in data collection, data preparation, model development, model evaluation, model post-processing, and model deployment (Baker & Hawn, 2022).

Tenth, once it becomes possible to identify what is wrong with a model through the transparency of its implementation, it may be followed by the accountability to address the harmful consequences of these algorithmic systems (Munn, 2022). This would entail a different discussion knowing that companies may address the harms they cause and be enforced by law to redesign their models.

10. References

- Abdullah, M., Madain, A., & Jararweh, Y. (2022). ChatGPT: Fundamentals, Applications and Social Impacts. *2022 Ninth International Conference on Social Networks Analysis, Management and Security (SNAMS)*, 1–8.
<https://doi.org/10.1109/SNAMS58071.2022.10062688>
- Alba, D. (2022, December 8). OpenAI Chatbot Spits Out Biased Musings, Despite Guardrails. *Bloomberg.Com*.
<https://www.bloomberg.com/news/newsletters/2022-12-08/chatgpt-open-ai-s-chatbot-is-spitting-out-biased-sexist-results>
- Aler Tubella, A., Barsotti, F., Koçer, R. G., & Mendez, J. A. (2022). Ethical implications of fairness interventions: What might be hidden behind engineering choices? *Ethics and Information Technology*, 24(1), 12.
- Alur, R., & Henzinger, T. A. (1998). Finitary fairness. *ACM Transactions on Programming Languages and Systems*, 20(6), 1171–1194.
<https://doi.org/10.1145/295656.295659>
- Amatriain, X. (2023). *Transformer models: An introduction and catalog* (arXiv:2302.07730). arXiv. <http://arxiv.org/abs/2302.07730>
- Apt, K. R., Francez, N., & Katz, S. (1988). Appraising fairness in languages for distributed programming. *Distributed Computing*, 2(4), 226–241.
<https://doi.org/10.1007/BF01872848>
- Ashokan, A., & Haas, C. (2021). Fairness metrics and bias mitigation strategies for rating predictions. *Information Processing & Management*, 58(5), 102646.
<https://doi.org/10.1016/j.ipm.2021.102646>
- Baharloo, A. (2013). Test Fairness in Traditional and Dynamic Assessment. *Theory and Practice in Language Studies*, 3(10), 1930–1938.
<https://doi.org/10.4304/tpls.3.10.1930-1938>
- Baker, R. S., & Hawn, A. (2022). Algorithmic Bias in Education. *International Journal of Artificial Intelligence in Education*, 32(4), 1052–1092.
<https://doi.org/10.1007/s40593-021-00285-9>
- Barocas, S., Crawford, K., Shapiro, A., & Wallach, H. (2017). The problem with bias: From allocative to representational harms in machine learning. Special Interest Group for Computing. *Information and Society (SIGCIS)*, 2.

- Baron, R. A. (2017). Income inequality. *Journal of Entrepreneurship and Public Policy*, 6(1), 2–10. <https://doi.org/10.1108/JEPP-07-2016-0028>
- BBC News. (2015, July 1). Google apologises for Photos app's racist blunder. *BBC News*. <https://www.bbc.com/news/technology-33347866>
- Berry, R. A. W. (2008). Novice teachers' conceptions of fairness in inclusion classrooms. *Teaching and Teacher Education*, 24(5), 1149–1159.
- Biddle, S. (2022, December 8). *The Internet's New Favorite AI Proposes Torturing Iranians and Surveilling Mosques*. The Intercept. <https://theintercept.com/2022/12/08/openai-chatgpt-ai-bias-ethics/>
- Binns, R. (2018). Fairness in Machine Learning: Lessons from Political Philosophy. *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, 149–159. <https://proceedings.mlr.press/v81/binns18a.html>
- Bisio, I., & Marchese, M. (2014). The concept of fairness: Definitions and use in bandwidth allocation applied to satellite environment. *IEEE Aerospace and Electronic Systems Magazine*, 29(3), 8–14.
- Brandao, M., Jirotko, M., Webb, H., & Luff, P. (2020). Fair navigation planning: A resource for characterizing and designing fairness in mobile robots. *Artificial Intelligence*, 282, 103259.
- Brockman, G., Atty, E., Georges, E., Jang, J., Kilpatrick, L., Lim, R., Miller, L., & Pokrass, M. (2023, January 3). *Introducing ChatGPT and Whisper APIs*. <https://openai.com/blog/introducing-chatgpt-and-whisper-apis>
- Buchanan, R. (1992). Wicked problems in design thinking. *Design Issues*, 8(2), 5–21.
- Card, D., & Smith, N. A. (2020). On Consequentialism and Fairness. *Frontiers in Artificial Intelligence*, 3. <https://www.frontiersin.org/articles/10.3389/frai.2020.00034>
- Castelnovo, A., Crupi, R., Greco, G., Regoli, D., Penco, I. G., & Cosentini, A. C. (2022). A clarification of the nuances in the fairness metrics landscape. *Scientific Reports*, 12(1), Article 1. <https://doi.org/10.1038/s41598-022-07939-1>
- Cellan-Jones, R. (2017, October 24). Facebook's News Feed experiment panics publishers. *BBC News*. <https://www.bbc.com/news/technology-41733119>
- Chan, S., & Zukerman, M. (2002). Is max-min fairness achievable in the presence of insubordinate users? *IEEE Communications Letters*, 6(3), 120–122. <https://doi.org/10.1109/4234.991152>
- Chen, J. S.-C., Cidon, I., & Ofek, Y. (1993). A local fairness algorithm for gigabit

- LAN's/MAN's with spatial reuse. *IEEE Journal on Selected Areas in Communications*, 11(8), 1183–1192. <https://doi.org/10.1109/49.245907>
- Chew, J. P., & Gupta, A. K. (2000). *Using dynamic weights for improving fairness in the ATM ABR service*. 372–377.
- Chhabra, A., Masalkovaitè, K., & Mohapatra, P. (2021). An Overview of Fairness in Clustering. *IEEE Access*, 9, 130698–130720. <https://doi.org/10.1109/ACCESS.2021.3114099>
- Christiano, P., Leike, J., Brown, T. B., Martic, M., Legg, S., & Amodei, D. (2023). *Deep reinforcement learning from human preferences* (arXiv:1706.03741). arXiv. <http://arxiv.org/abs/1706.03741>
- Cleary, T. A. (1968). Test Bias: Prediction of Grades of Negro and White Students in Integrated Colleges. *Journal of Educational Measurement*, 5(2), 115–124.
- Cole, N. S. (1973). Bias in Selection. *Journal of Educational Measurement*, 10(4), 237–255.
- Cole, N. S., & Zieky, M. J. (2001). The new faces of fairness. *Journal of Educational Measurement*, 38(4), 369–382.
- Conklin, J. (2006). *Wicked problems & social complexity* (Vol. 11). CogNexus Institute Napa, USA.
- Coyne, R. (2005). Wicked problems revisited. *Design Studies*, 26(1), 5–17.
- Crawford, K. (Director). (2017, December 10). *The Trouble with Bias*. https://www.youtube.com/watch?v=fMym_BKWQzk
- Crocker, L. (2003). Teaching for the Test: Validity, Fairness, and Moral Action. *Educational Measurement: Issues and Practice*, 22(3), 5–11. <https://doi.org/10.1111/j.1745-3992.2003.tb00132.x>
- Darlington, R. B. (1971). Another Look at “Cultural Fairness.” *Journal of Educational Measurement*, 8(2), 71–82.
- Das, K. (2021). *The Paradigm of Justice: A Contemporary Debate between John Rawls and Amartya Sen*. Routledge.
- de Oliveira, L. F., da Silva Gomes, A., Enes, Y., Castelo Branco, T. V., Pires, R. P., Bolzon, A., & Demo, G. (2022). Path and future of artificial intelligence in the field of justice: A systematic literature review and a research agenda. *SN Social Sciences*, 2(9), 180.
- Dehouche, N. (2021). Implicit Stereotypes in Pre-Trained Classifiers. *IEEE ACCESS*, 9, 167936–167947. <https://doi.org/10.1109/ACCESS.2021.3136898>
- Deutsch, M. (1985). *Distributive justice: A social-psychological perspective*.

- Duran-Rodas, D., Villeneuve, D., Pereira, F. C., & Wulforst, G. (2020). How fair is the allocation of bike-sharing infrastructure? Framework for a qualitative and quantitative spatial fairness assessment. *Transportation Research Part A: Policy and Practice*, 140, 299–319. <https://doi.org/10.1016/j.tra.2020.08.007>
- Durden, T. (2023, March 23). *Google's New Bard AI Is Riddled With Political Bias*. ZeroHedge. <https://www.zerohedge.com/technology/googles-new-bard-ai-riddled-political-bias>
- Esmer, S. (2021). *Amartya Sen's capability approach and its relation with John Rawls 'justice as fairness*.
- Fei, Z., & Yang, M. (2005). Intra-Session Fairness in Multicast Communications. *Telecommunication Systems*, 29(4), 235–255. <https://doi.org/10.1007/s11235-005-3268-9>
- Ferrara, E. (2023). Should chatgpt be biased? Challenges and risks of bias in large language models. *ArXiv Preprint ArXiv:2304.03738*.
- Flaughter, R. L. (1974). The New Definitions of Test Fairness in Selection: Developments and Implications. *Educational Researcher*, 3(9), 13–16. <https://doi.org/10.2307/1174915>
- Fletcher, R. R., Nakeshimana, A., & Olubeko, O. (2021). Addressing Fairness, Bias, and Appropriate Use of Artificial Intelligence and Machine Learning in Global Health. *Frontiers in Artificial Intelligence*, 3. <https://www.frontiersin.org/articles/10.3389/frai.2020.561802>
- Friedler, S. A., Scheidegger, C., & Venkatasubramanian, S. (2021). The (Im)possibility of fairness: Different value systems require different mechanisms for fair decision making. *Communications of the ACM*, 64(4), 136–143. <https://doi.org/10.1145/3433949>
- Gabriel, I. (2022). Toward a theory of justice for artificial intelligence. *Daedalus*, 151(2), 218–231.
- Gao, R., & Shah, C. (2020). Toward creating a fairer ranking in search engine results. *Information Processing & Management*, 57(1), 102138.
- Garay, J., MacKenzie, P., Prabhakaran, M., & Yang, K. (2006). *Resource fairness and composability of cryptographic protocols*. 404–428.
- Garske, V. (2023, May 20). *Transformers—Timeline of Transformer Models*.

- <https://ai.v-gar.de/ml/transformer/timeline/attention.html>
- Ghojogh, B., & Ghodsi, A. (2020). *Attention mechanism, transformers, BERT, and GPT: Tutorial and survey*.
- Giovanola, B., & Tiribelli, S. (2022). Weapons of moral construction? On the value of fairness in algorithmic decision-making. *Ethics and Information Technology*, 24(1), 3.
- Goldman, B., & Cropanzano, R. (2015). "Justice" and "fairness" are not the same thing. *Journal of Organizational Behavior*, 36(2), 313–318.
- Green, B. (2022). Escaping the Impossibility of Fairness: From Formal to Substantive Algorithmic Fairness. *Philosophy & Technology*, 35(4), 90.
<https://doi.org/10.1007/s13347-022-00584-6>
- Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C. C. T., Del Giorno, A., Gopi, S., Javaheripi, M., Kauffmann, P., de Rosa, G., Saarikivi, O., Salim, A., Shah, S., Behl, H. S., Wang, X., Bubeck, S., Eldan, R., Kalai, A. T., Lee, Y. T., & Li, Y. (2023). *Textbooks Are All You Need* (arXiv:2306.11644). arXiv. <http://arxiv.org/abs/2306.11644>
- Gupta, S., & Kamble, V. (2021). Individual fairness in hindsight. *The Journal of Machine Learning Research*, 22(1), 144:6386-144:6420.
- Hall, M. R. (2021). The living wage gap—A quantitative measure of poverty in global supply chains. *The International Journal of Life Cycle Assessment*, 26(9), 1867–1877. <https://doi.org/10.1007/s11367-021-01945-7>
- Hamermesh, D. S., & Schmidt, P. (2003). The determinants of econometric society fellows elections. *Econometrica*, 399–407.
- Häyry, M. (2018). Doctrines and Dimensions of Justice: Their Historical Backgrounds and Ideological Underpinnings. *Cambridge Quarterly of Healthcare Ethics*, 27(2), 188–216. <https://doi.org/10.1017/S096318011700055X>
- He, P., Gao, J., & Chen, W. (2023). *DeBERTaV3: Improving DeBERTa using ELECTRA-Style Pre-Training with Gradient-Disentangled Embedding Sharing* (arXiv:2111.09543). arXiv. <https://doi.org/10.48550/arXiv.2111.09543>
- Howard, R. J., Tallontire, A. M., Stringer, L. C., & Marchant, R. A. (2016). Which "fairness", for whom, and why? An empirical analysis of plural notions of fairness in Fairtrade Carbon Projects, using Q methodology. *Environmental Science & Policy*, 56, 100–109. <https://doi.org/10.1016/j.envsci.2015.11.009>
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021).

- LoRA: Low-Rank Adaptation of Large Language Models* (arXiv:2106.09685). arXiv.
<https://doi.org/10.48550/arXiv.2106.09685>
- Jensen, A. (1969). How Much Can We Boost IQ and Scholastic Achievement. *Harvard Educational Review*, 39(1), 1–123.
<https://doi.org/10.17763/haer.39.1.l3u15956627424k7>
- Johnson, M. S., Liu, X., & McCaffrey, D. F. (2022). Psychometric Methods to Evaluate Measurement and Algorithmic Bias in Automated Scoring. *Journal of Educational Measurement*, 59(3), 338–361. <https://doi.org/10.1111/jedm.12335>
- Khashabi, D., Min, S., Khot, T., Sabharwal, A., Tafjord, O., Clark, P., & Hajishirzi, H. (2020). UNIFIEDQA: Crossing Format Boundaries with a Single QA System. *Findings of the Association for Computational Linguistics: EMNLP 2020*, 1896–1907.
<https://doi.org/10.18653/v1/2020.findings-emnlp.171>
- Kirk, H. R., Jun, Y., Volpin, F., Iqbal, H., Benussi, E., Dreyer, F., Shtedritski, A., & Asano, Y. (2021). Bias Out-of-the-Box: An Empirical Analysis of Intersectional Occupational Biases in Popular Generative Language Models. *Advances in Neural Information Processing Systems*, 34, 2611–2624.
<https://proceedings.neurips.cc/paper/2021/hash/1531beb762df4029513ebf9295e0d34f-Abstract.html>
- Kleinberg, J., Mullainathan, S., & Raghavan, M. (2016). Inherent trade-offs in the fair determination of risk scores. *ArXiv Preprint ArXiv:1609.05807*.
- Krysiak, F. C., & Krysiak, D. (2006). Sustainability with Uncertain Future Preferences. *Environmental and Resource Economics*, 33(4), 511–531.
<https://doi.org/10.1007/s10640-005-0004-6>
- Kwiatkowska, M. Z. (1989). Survey of fairness notions. *Information and Software Technology*, 31(7), 371–386. [https://doi.org/10.1016/0950-5849\(89\)90159-6](https://doi.org/10.1016/0950-5849(89)90159-6)
- Ledvinka, J. (1979). The Statistical Definition of Fairness in the Federal Selection Guidelines and Its Implications for Minority Employment. *Personnel Psychology*, 32(3), 551–562. <https://doi.org/10.1111/j.1744-6570.1979.tb02153.x>
- Lee, M. S. A., & Floridi, L. (2021). Algorithmic Fairness in Mortgage Lending: From Absolute Conditions to Relational Trade-offs. *Minds and Machines*, 31(1), 165–191. <https://doi.org/10.1007/s11023-020-09529-4>
- Lee, N. T., Resnick, P., & Barton, G. (2019). Algorithmic bias detection and mitigation: Best practices and policies to reduce consumer harms. *Brookings Institute*:

Washington, DC, USA, 2.

- Li, J.-S., Liu, C.-G., & Huang, C.-Y. (2007). Achieving multipoint-to-multipoint fairness with RCNWA. *Journal of Systems Architecture*, 53(7), 437–452.
- Li, L., Wang, L., Chen, J., Wang, R., & Zhang, Z. (2014). Fairness Analysis for Multiparty Nonrepudiation Protocols Based on Improved Strand Space. *Discrete Dynamics in Nature and Society*, 2014, 1–7. <https://doi.org/10.1155/2014/904717>
- Li, T., Khashabi, D., Khot, T., Sabharwal, A., & Srikumar, V. (2020). UNQOVERing Stereotyping Biases via Underspecified Questions. *Findings of the Association for Computational Linguistics: EMNLP 2020*, 3475–3489. <https://doi.org/10.18653/v1/2020.findings-emnlp.311>
- Li, Y., Huang, H., Guo, X., & Yuan, Y. (2021). An Empirical Study on Group Fairness Metrics of Judicial Data. *IEEE ACCESS*, 9, 149043–149049. <https://doi.org/10.1109/ACCESS.2021.3122443>
- Li, Y., & Zhang, Y. (2023). *Fairness of ChatGPT* (arXiv:2305.18569). arXiv. <http://arxiv.org/abs/2305.18569>
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). *RoBERTa: A Robustly Optimized BERT Pretraining Approach* (arXiv:1907.11692). arXiv. <https://doi.org/10.48550/arXiv.1907.11692>
- Lundgard, A. (2020). Measuring Justice in Machine Learning. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 680. <https://doi.org/10.1145/3351095.3372838>
- MacKinnon, C. A. (2011). Substantive Equality: A Perspective. *Minnesota Law Review*, 96, 1.
- Maffulli, S. (2023, July 20). Meta’s LLaMa 2 license is not Open Source. *Voices of Open Source*. <https://blog.opensource.org/metast-llama-2-license-is-not-open-source/>
- Markowitch, O., Gollmann, D., & Kremer, S. (2003). On Fairness in Exchange Protocols. In P. J. Lee & C. H. Lim (Eds.), *Information Security and Cryptology—ICISC 2002* (Vol. 2587, pp. 451–465). Springer Berlin Heidelberg. https://doi.org/10.1007/3-540-36552-4_31
- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A Survey on Bias and Fairness in Machine Learning. *ACM Comput. Surv.*, 54(6). <https://doi.org/10.1145/3457607>
- Moret, F., Pinson, P., & Papakonstantinou, A. (2020). Heterogeneous risk preferences in

- community-based electricity markets. *European Journal of Operational Research*, 287(1), 36–48. <https://doi.org/10.1016/j.ejor.2020.04.034>
- Munn, L. (2022). The uselessness of AI ethics. *AI and Ethics*.
<https://doi.org/10.1007/s43681-022-00209-w>
- Munn, L. (2023). The five tests: Designing and evaluating AI according to indigenous Māori principles. *AI & SOCIETY*. <https://doi.org/10.1007/s00146-023-01636-x>
- Nguyen, U. T., & Katzela, I. (2000). *A flexible multipoint-to-point traffic control algorithm for ABR services in ATM networks*. 185–194.
- Nicosia, G., Pacifici, A., & Pferschy, U. (2017). Price of Fairness for Allocating a Bounded Resource. *European Journal of Operational Research*, 257(3), 933–943.
<https://doi.org/10.1016/j.ejor.2016.08.013>
- Non-deterministic*. (2023, July 12).
<https://dictionary.cambridge.org/dictionary/english/non-deterministic>
- Nussbaum, M. C. (2011). *Creating Capabilities: The Human Development Approach*. Harvard University Press. <https://doi.org/10.2307/j.ctt2jbt31>
- OpenAI. (2017, June 13). *Learning from human preferences*.
<https://openai.com/research/learning-from-human-preferences>
- OpenAI. (2021, June 10). *Improving language model behavior by training on a curated dataset*. <https://openai.com/research/improving-language-model-behavior>
- OpenAI. (2022, November 30). *Introducing ChatGPT*. <https://openai.com/blog/chatgpt>
- OpenAI. (2023a). *GPT - OpenAI API*.
<https://platform.openai.com/docs/guides/gpt/chat-completions-api>
- OpenAI. (2023b, 02). *Introducing ChatGPT Plus*. <https://openai.com/blog/chatgpt-plus>
- OpenAI. (2023c, February 16). *How should AI systems behave, and who should decide?*
<https://openai.com/blog/how-should-ai-systems-behave>
- Österberg, P., & Zhang, T. (2011). Multicast-Favourable Max-Min Fairness–The Definition and How to Comply. *International Journal of Computers and Applications*, 33(1), 1–8.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askeel, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022). *Training language models to follow instructions with human feedback* (arXiv:2203.02155). arXiv.
<http://arxiv.org/abs/2203.02155>

- Parrish, A., Chen, A., Nangia, N., Padmakumar, V., Phang, J., Thompson, J., Htut, P. M., & Bowman, S. R. (2022). *BBQ: A Hand-Built Bias Benchmark for Question Answering* (arXiv:2110.08193). arXiv. <http://arxiv.org/abs/2110.08193>
- Pasquali, L. (2009). Psychometrics. *Revista Da Escola de Enfermagem Da USP*, 43, 992–999. <https://doi.org/10.1590/S0080-62342009000500002>
- Paulus, J. K., & Kent, D. M. (2020). Predictably unequal: Understanding and addressing concerns that algorithmic clinical prediction may increase health disparities. *Npj Digital Medicine*, 3(1), Article 1. <https://doi.org/10.1038/s41746-020-0304-9>
- Persico, N. (2002). Racial Profiling, Fairness, and Effectiveness of Policing. *The American Economic Review*, 92(5), 1472–1497.
- Phi, M. (Director). (2020, April 28). *Illustrated Guide to Transformers Neural Network: A step by step explanation*. <https://www.youtube.com/watch?v=4Bdc55j80l8>
- Piccininni, M. (2022). Counterfactual fairness: The case study of a food delivery platform’s reputational-ranking algorithm. *Frontiers in Psychology*, 13.
- Pichai. (2023, February 6). *An important next step on our AI journey*. Google. <https://blog.google/technology/ai/bard-google-ai-search-updates/>
- Queille, J. P., & Sifakis, J. (1983). Fairness and related properties in transition systems—A temporal logic to deal with fairness. *Acta Informatica*, 19(3), 195–220. <https://doi.org/10.1007/BF00265555>
- Rasooli, A., Zandi, H., & DeLuca, C. (2019). Conceptualising fairness in classroom assessment: Exploring the value of organisational justice theory. *Assessment in Education: Principles, Policy & Practice*, 26(5), 584–611.
- Rawls, J. (1971). *A theory of justice*. Cambridge (Mass.).
- Rawls, J. (1978). *The basic structure as subject*. Springer.
- Rawls, J. (1996). *Political Liberalism, the John Dewy Essays in Philosophy, Number Four*.
- Ray, P. P. (2023). ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems*, 3, 121–154. <https://doi.org/10.1016/j.iotcps.2023.04.003>
- Reinl, J. (2023, March 27). *Google’s chatbot denies bias, but promotes transgenderism and veganism*. Mail Online. <https://www.dailymail.co.uk/news/article-11908383/Googles-chatbot-denies-bias-promotes-transgenderism-Joe-Biden-veganism.html>

- Rittel, H. W., & Webber, M. M. (1973). Dilemmas in a general theory of planning. *Policy Sciences*, 4(2), 155–169.
- Robeyns, I. (2006). The capability approach in practice. *Journal of Political Philosophy*, 14(3), 351–376.
- Robeyns, I. (2009). Justice as fairness and the capability approach. *Arguments for a Better World: Essays in Honor of Amartya Sen*, 1, 397–413.
- Rubenstein, D., Kurose, J., & Towsley, D. (2002). The impact of multicast layering on network fairness. *IEEE/ACM Transactions on Networking*, 10(2), 169–182.
<https://doi.org/10.1109/90.993299>
- Saito, M. (2003). Amartya Sen’s capability approach to education: A critical exploration. *Journal of Philosophy of Education*, 37(1), 17–33.
- Salimi, B., Howe, B., & Suci, D. (2020). Database Repair Meets Algorithmic Fairness. *SIGMOD Rec.*, 49(1), 34–41. <https://doi.org/10.1145/3422648.3422657>
- Saxena, N. A., Huang, K., DeFilippis, E., Radanovic, G., Parkes, D. C., & Liu, Y. (2020). How do fairness definitions fare? Testing public attitudes towards three algorithmic definitions of fairness in loan allocations. *ARTIFICIAL INTELLIGENCE*, 283.
<https://doi.org/10.1016/j.artint.2020.103238>
- Scantamburlo, T. (2021). Non-empirical problems in fair machine learning. *ETHICS AND INFORMATION TECHNOLOGY*, 23(4), 703–712.
<https://doi.org/10.1007/s10676-021-09608-9>
- Scutari, M., Panero, F., & Proissl, M. (2022). Achieving fairness with a simple ridge penalty. *Statistics and Computing*, 32(5), 77.
<https://doi.org/10.1007/s11222-022-10143-w>
- Sen, A. (1980). *Equality of what?*
- Sen, A. (1988). *The standard of living*. Cambridge University Press.
- Sen, A. (1990). Justice: Means versus Freedoms. *Philosophy & Public Affairs*, 19(2), 111–121.
- Sen, A. (2008). The idea of justice. *Journal of Human Development*, 9(3), 331–342.
- Sorkin, A. R., Mattu, R., Kessler, S., Merced, M. J. de la, Hirsch, L., & Livni, E. (2023, February 23). China Chases ChatGPT’s Success. *The New York Times*.
<https://www.nytimes.com/2023/02/23/business/dealbook/china-ai-chatgpt.html>
- Steffy, B. D., & Ledvinka, J. (1989). The long-range impact of five definitions of “fair”

- employee selection on black employment and employee productivity. *Organizational Behavior and Human Decision Processes*, 44(2), 297–324.
[https://doi.org/10.1016/0749-5978\(89\)90029-0](https://doi.org/10.1016/0749-5978(89)90029-0)
- Suresh, H., & Guttag, J. V. (2021). A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. *Equity and Access in Algorithms, Mechanisms, and Optimization*, 1–9. <https://doi.org/10.1145/3465416.3483305>
- Szepannek, G., & Lübke, K. (2021). Facing the challenges of developing fair risk scoring models. *Frontiers in Artificial Intelligence*, 4, 681915.
- Topal, M. O., Bas, A., & van Heerden, I. (2021). *Exploring Transformers in Natural Language Generation: GPT, BERT, and XLNet* (arXiv:2102.08036). arXiv.
<http://arxiv.org/abs/2102.08036>
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., & Lample, G. (2023). *LLaMA: Open and Efficient Foundation Language Models* (arXiv:2302.13971). arXiv. <http://arxiv.org/abs/2302.13971>
- UBS Editorial Team. (2023, February 23). *Let's chat about ChatGPT*. Global.
<https://www.ubs.com/global/en/wealth-management/our-approach/marketnews/article.1585717.html>
- van Nood, R., & Yeomans, C. (2021). Fairness as Equal Concession: Critical Remarks on Fair AI. *Science and Engineering Ethics*, 27(6), 73.
<https://doi.org/10.1007/s11948-021-00348-z>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is All you Need. *Advances in Neural Information Processing Systems*, 30.
<https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>
- Verma, S., & Rubin, J. (2018). Fairness definitions explained. *Proceedings of the International Workshop on Software Fairness*, 1–7.
<https://doi.org/10.1145/3194770.3194776>
- Vock, I. (2022, December 9). ChatGPT proves that AI still has a racism problem. *New Statesman*.
<https://www.newstatesman.com/science-tech/2022/12/chatgpt-shows-ai-racism-problem>

- Wabenhorst, A. (2003). Stepwise development of fair distributed systems. *Acta Informatica*, 39, 233–271.
- Walker, M. (2005). Amartya Sen’s capability approach and education. *Educational Action Research*, 13(1), 103–110.
- Wang, E. J. (2023). *Alpaca-LoRA* [Jupyter Notebook].
<https://github.com/tloen/alpaca-lora> (Original work published 2023)
- Wang, P., Jiang, H., Zhuang, W., & Poor, H. V. (2008). Redefinition of max-min fairness in multi-hop wireless networks. *IEEE Transactions on Wireless Communications*, 7(12), 4786–4791.
- Weidinger, L., McKee, K. R., Everett, R., Huang, S., Zhu, T. O., Chadwick, M. J., Summerfield, C., & Gabriel, I. (2023). Using the Veil of Ignorance to align AI systems with principles of justice. *Proceedings of the National Academy of Sciences*, 120(18), e2213709120.
- Woodruff, A., Fox, S. E., Rousso-Schindler, S., & Warshaw, J. (2018). A Qualitative Exploration of Perceptions of Algorithmic Fairness. *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, 1–14.
<https://doi.org/10.1145/3173574.3174230>
- Xia, B., Yin, J., Xu, J., & Li, Y. (2019). WE-Rec: A fairness-aware reciprocal recommendation based on Walrasian equilibrium. *Knowledge-Based Systems*, 182, 104857. <https://doi.org/10.1016/j.knosys.2019.07.028>
- Zehlike, M., Hacker, P., & Wiedemann, E. (2020). Matching code and law: Achieving algorithmic fairness with optimal transport. *Data Mining and Knowledge Discovery*, 34(1), 163–200.
- Zehlike, M., Sühr, T., Baeza-Yates, R., Bonchi, F., Castillo, C., & Hajian, S. (2022). Fair Top-k Ranking with multiple protected groups. *Information Processing & Management*, 59(1), 102707. <https://doi.org/10.1016/j.ipm.2021.102707>
- Zhang, J., Bao, K., Zhang, Y., Wang, W., Feng, F., & He, X. (2023). *Is ChatGPT Fair for Recommendation? Evaluating Fairness in Large Language Model Recommendation* (arXiv:2305.07609). arXiv. <http://arxiv.org/abs/2305.07609>
- Zhuo, T. Y., Huang, Y., Chen, C., & Xing, Z. (2023). *Red teaming ChatGPT via Jailbreaking: Bias, Robustness, Reliability and Toxicity* (arXiv:2301.12867). arXiv.
<http://arxiv.org/abs/2301.12867>
- Zollers, N. J., Albert, L. R., & Cochran-Smith, M. (2000). In Pursuit of Social Justice:

- Collaborative Research and Practice in Teacher Education. *Action in Teacher Education*, 22(2), 1–14. <https://doi.org/10.1080/01626620.2000.10463000>
- Zukerman, M., & Chan, S. (1993). Fairness in ATM networks. *Computer Networks and ISDN Systems*, 26(1), 109–117.
- Zwick, R. (2019). Fairness in measurement and selection: Statistical, philosophical, and public perspectives. *Educational Measurement: Issues and Practice*, 38(4), 34–41.

11. Appendix

11.1. Appendix A: Transformer model architecture

Vaswani et al. (2017), in their paper *Attention is All You Need*, introduce the transformer model as an instance of the encoder-decoder architecture that uses attention mechanisms for sequence-to-sequence modelling. By then, recurrent neural networks (RNN) and long short-term memory (LSTM) networks were established as the state-of-the-art approaches for sequence modelling. However, these models had significant difficulty with long sequences as the probability of maintaining context between distant words in a sentence decreased exponentially (Topal et al., 2021). The key advantage of transformer models is then that they rely entirely on the attention mechanism to create global dependencies between inputs and outputs (Amatriain, 2023).

Since the initial paper, there has been a large increase in the development of new models and methods using a similar architecture. However, despite sharing the same architecture, many of these new models do not have self-explanatory names (Amatriain, 2023). As a consequence, some authors have developed extensive lists of models that use the transformer architecture and the quick development of new models requires further updates (Garske, 2023).

Encoder and decoder architecture

The input feeds a word embedding layer that can be seen as a lookup table to learn a numerical representation of a word in vectors (Phi, 2020). Since the model has no recurrence and no convolution and the order of words is important for the meaning of sentences, these vectors are used to account for the position to each input word embedding (Ghojogh & Ghodsi, 2020). In order to use the sequence of the input, Vaswani et al. (2017) explains that positional information should be injected into the tokens of the sequence. This process is called positional encoding in which, in transformer models, a vector of the same size of the input embeddings d_{model} is generated for every position of the input and every dimension within the embedding vector is calculated with sine and cosine functions. The sine and cosine functions are

used for even and odd dimensions respectively. This positional encoding is later added to the embedding vectors to give information of the positions of each token.

Once this process is done, this information is given to the encoder of the model. The encoder consists of 2 sublayers, the first one is a multi-head self-attention module and the second is feed-forward network (Vaswani et al., 2017). The multi-head self-attention sublayer allows to associate the importance of each word in the input to the other words in the input (Phi, 2020). Phi (2020) explains that to achieve this, the inputs from the previous process are fed to 3 connected layers to map a query to a set of key-value pairs where query, keys, and values are all vectors. The query and key vectors follow a dot product multiplication to create a score matrix which gives a numerical value to the attention a word should put into another word. The higher the number the more the focus. The attention scores are scaled down by applying the softmax function in order to create a better differentiation between the words that are relevant and the rest. Later, the result of this process, a matrix with the attention weights, will be multiplied by the value matrix, containing the word embeddings for every token. The result of this multiplication is fed to a linear layer to process. To make this a multi-headed process, the three vectors need to be split into adding vectors that go through the self-attention process described earlier, each one being considered a head. After the multi-headed attention process is complete, the output is added to the original input and fed to a normalising and feed forward layer to obtain a higher representation of the input. This process is done in order to encoding the input to a continuous representation with attention information.

The decoder is used to generate word sequences (Phi, 2020). Vaswani et al. (2017) explain that the encoder adds a third sub-layer that performs multi-headed attention with the output from the encoder module. Moreover, the decode module modifies the self-attention sublayer with masking in order to prevent positions from paying attention to subsequent positions and return output embeddings with a position offset. Phi explains that the decoder is auto regressive by taking previous decoder outputs along with the encoder outputs as inputs. In the encoder, the input goes through positional encoding to get positional embeddings that are later fed to the first multi-headed attention sublayer. However, this sublayer is slightly different to the one of the encoder, as the decoder generates de sequence word by word, it needs to be prevented from

conditioning from future tokens meaning that a word should have only access to itself and the words generated before. This process is called masking, the sum of a mask matrix of the same size to the attention scores matrix. The result of this process is a masked vector with information on what attention to pay to the inputs of the decoder. This information is fed to the second sublayer of multi-headed attention. The output of the encoder are keys and values while the previous layer output acts as the query. Once such a process is done the output goes through a final linear classifier with the amount of classes equal to the amount of words in the vocabulary. The classified output is fed to a softmax that defines a probability to each class (word in the vocabulary) and takes the index of the highest probability to define the predicted word. This process happens until the end token is predicted.

Alignment

According to Ouyang et al. (2022), making large language models does not fix the problems of returning untruthful, toxic or not helpful outputs. The authors argue that the objective for many recent language models of predicting the next token from the internet, is different to the one of following instructions helpfully and safely. In 2017, OpenAI in collaboration with DeepMind developed an algorithm that infers what is the desired output based on selecting one of two outputs in a scalable manner (OpenAI, 2017). The approach of OpenAI is to learn a reward function from human feedback and then optimise that reward function (Christiano et al., 2023). The approach requires selecting one from multiple options instead of giving a numerical score.

In GPT models, RLHF has been used by collecting a dataset of human-labelled comparisons from outputs of the models (Ouyang et al., 2022). This dataset is later used to train the reward model that is applied as a reward function to fine-tune the learning baseline using the Proximal policy optimization algorithm. According to Ouyang et al. (2022), the cost of increasing the alignment of these models using RLHF is modest compared to the pretraining cost. At the same time, RLHF makes language models more helpful up to 100 times the model size increase. Furthermore, models fine-tuned with this framework seem to generalise following instructions to settings in which they were not trained.