

Study on the impact of artificial intelligence systems, their potential for promoting equality, including gender equality, and the risks they may cause in relation to non-discrimination



GENDER EQUALITY COMMISSION (GEC)
AND THE STEERING COMMITTEE
ON ANTI-DISCRIMINATION,
DIVERSITY AND INCLUSION (CDADI)

COUNCIL OF EUROPE



CONSEIL DE L'EUROPE

Study on the impact of artificial intelligence systems, their potential for promoting equality, including gender equality, and the risks they may cause in relation to non-discrimination

**GENDER EQUALITY COMMISSION (GEC)
AND THE STEERING COMMITTEE ON
ANTI-DISCRIMINATION, DIVERSITY
AND INCLUSION (CDADI)**

**Prepared by Ivana Bartoletti
and Raphaële Xenidis**

French edition:

Étude sur l'impact des systèmes d'intelligence artificielle, leur potentiel de promotion de l'égalité, y compris l'égalité de genre, et les risques qu'ils peuvent entraîner en matière de non-discrimination

The opinions expressed in this work are the responsibility of the author(s) and do not necessarily reflect the official policy of the Council of Europe.

The reproduction of extracts (up to 500 words) is authorised, except for commercial purposes as long as the integrity of the text is preserved, the excerpt is not used out of context, does not provide incomplete information or does not otherwise mislead the reader as to the nature, scope or content of the text. The source text must always be acknowledged as follows "© Council of Europe, year of the publication". All other requests concerning the reproduction/translation of all or part of the document, should be addressed to the Directorate of Communications, Council of Europe (F-67075 Strasbourg Cedex or publishing@coe.int). All other correspondence concerning this document should be addressed to the Anti-Discrimination Department of the Council of Europe.

Council of Europe
F-67075 Strasbourg Cedex France
E-mail: cdadi@coe.int

Cover design and layout: Documents and Publications Production Department (SPDP), Council of Europe
Cover photo: IStock

This publication has not been copy-edited by the SPDP Editorial Unit to correct typographical and grammatical errors.

© Council of Europe, August 2023
Printed at the Council of Europe

Prepared by:

Ivana Bartoletti,

Global Chief Privacy Officer at Wipro, Visiting Policy Fellow at the Oxford Internet Institute, University of Oxford and co-founded the Women Leading in AI Network.

Raphaële Xenidis,

Assistant Professor in European law, Sciences Po, Law School.

Contents

EXECUTIVE SUMMARY	5
INTRODUCTION: THE CONTEXT	9
SECTION 1	13
UNPACKING ‘MACHINE BIAS’: HOW CAN ALGORITHMIC TECHNOLOGIES LEAD TO DISCRIMINATION?	13
1) What is AI?	13
2) What is algorithmic bias?	15
3) The discriminatory impact of AI: some concrete examples	19
4) What makes algorithmic discrimination different?	28
5) Addressing algorithmic discrimination: best practices and their limits	30
6) Representation and participation issues: The lack of diversity and inclusion in the AI industry	36
SECTION 2	41
THE LEGAL AND POLICY LANDSCAPE IN EUROPE: STRENGTHS AND SHORTCOMINGS	41
I. Discrimination and equality: legal and policy instruments and their limits	42
II. Privacy and data protection law: Fairness and accuracy	59
III. AI sectoral regulations: strengths and limits for promoting equality and addressing discrimination	62
SECTION 3	67
PROMOTING EQUALITY IN AND THROUGH THE USE OF AI: THE ROLE OF POSITIVE ACTION AND POSITIVE OBLIGATIONS	67
I. Revisiting existing rules in light of new power asymmetries	67
II. An obligation to promote equality in and through the use of algorithmic systems: the role of positive action and positive obligations	72
SECTION 4	79
RECOMMENDATIONS	79
Policy recommendations: Towards human-rights-based approach to AI	79

Executive summary

As the deployment of algorithmic systems and AI applications is growing in size and significance, algorithmic discrimination has become a matter of rising public concern. Regulatory responses are currently being devised across the globe, including in the European Union. The Council of Europe has initiated work on a legal framework for the development, design and application of artificial intelligence, based on the Council of Europe standards on human rights, democracy and the rule of law. Based on a “Feasibility Study on legal framework on AI design, development and application based on Council of Europe standards” published in 2020 as well as the “Possible elements of a legal framework on artificial intelligence, based on the Council of Europe’s standards on human rights, democracy and the rule of law”, the Committee on Artificial Intelligence (CAI), set up in 2022, is in the process of drafting a Framework Convention “on the development, design, and application of artificial intelligence systems based on the Council of Europe’s standards on human rights, democracy and the rule of law, and conducive to innovation”.¹ Such a legally binding instrument of the Council of Europe has the potential to foster a **human-rights-based approach** to the use of AI and algorithmic technologies in and beyond the international community of State Parties to the European Convention on Human Rights (ECHR). In this perspective, this study investigates the discriminatory risks of algorithmic technologies, the specific legal responses to algorithmic discrimination that can be offered by the Council of Europe, and the potential of these technologies for promoting equality, including gender equality. The Study is structured in three sections followed by recommendations, a summary of which is provided below. The first section unpacks issues of machine bias and reviews how algorithmic technologies can lead to discrimination. The second section probes the strengths and shortcomings of the legal framework that can be relied on to address algorithmic discrimination at the Council of Europe level. The third section investigates how positive action and positive obligations can be used to tackle algorithmic discrimination from its social roots to its manifestations in technological deployments in a transformative manner.

1. See Terms of reference of the Committee on Artificial Intelligence CM(2021)131 available at: <https://rm.coe.int/cai-terms-of-reference/1680a7b90b>.

Summary of recommendations

As highlighted in this Study, addressing the problem of algorithmic discrimination requires a multi-faceted response. This Study suggests that the Council of Europe should develop a robust human-rights-based approach to AI in the field of equality through the preparation of a specific Committee of Ministers Recommendation on AI, equality, including gender equality, and discrimination. This instrument should be drafted by an expert committee under the Gender Equality Commission (GEC) and the Steering Committee on Anti-Discrimination, Diversity and Inclusion (CDADI) and would build on the general human rights principles of equality, including gender equality and non-discrimination, including as they will appear in a future Framework Convention on artificial intelligence which is currently under preparation. This approach should include four complementary avenues for regulatory and policy intervention:

1. Prevention, transparency and accountability

Member states are encouraged to:

- ▶ expand the use of **positive action** measures to tackle algorithmic discrimination and to use the concept of positive obligations anchored in ECHR case law to create an obligation for providers and users to reasonably prevent algorithmic discrimination.
- ▶ introduce **mandatory discrimination risk and equality impact assessments** throughout the lifecycle of algorithmic systems according to their specific uses.
- ▶ consider how **certification mechanisms** could be used to ensure that biases have been mitigated and risks of discrimination eliminated as far as possible for well-defined uses.
- ▶ investigate the **relationship between accountability, transparency and trade secrets law** as it pertains to AI and the risks that it may pose to equality, including gender equality, and non-discrimination.
- ▶ consider establishing legal obligations for users of AI systems to **publish statistical data** that can allow interested parties to assess the discriminatory effects of a given system in the context of discrimination claims.
- ▶ introduce **mechanisms for transparency** with a view to allowing interested persons to assess potential discriminatory effects of a given system.
- ▶ consolidate prevention, transparency and accountability measures in a comprehensive **Action Plan on AI and Equality**.

2. Access to justice and legal redress mechanisms

Member states are encouraged to:

- ▶ facilitate access to justice by establishing **public supervision mechanisms** and developing **collective action routes** for redress of algorithmic discrimination.
- ▶ consider adjusting, complementing and reinforcing the effectiveness of **evidence rules** to create a **fairer, more balanced burden of proof**.
- ▶ encourage **co-operation between regulatory bodies and agencies**.
- ▶ investigate the **new forms of 'algorithmic' vulnerability** that emerge with the use of AI systems and consider **legal protection** against such vulnerability.
- ▶ make clear that the prohibition of discrimination in Art.14 ECHR covers **intersectional discrimination and discrimination by proxy**, two forms of discriminatory harms that algorithmic systems are most likely to generate.
- ▶ explore how **consumer protection law** could be used to **complement anti-discrimination law**, for instance by facilitating access to information, prohibiting certain features in algorithmic systems under the notion of abusive clauses, etc.

3. Diversity, inclusion, representation and participation

- ▶ Member states should **identify, support and actively enforce positive action measures**, including measures diversifying professional communities.
- ▶ **Positive obligations to promote equality** should provide the legal basis to ensure that AI and algorithmic systems are developed with equality promotion at their core.
- ▶ Positive obligations to promote equality could also translate into a requirement for companies of the AI sector to **develop and implement an equality strategy** covering the groups protected under Article 14 ECHR and Article 1 Protocol no. 12 ECHR.

4. Democratic participation, public awareness-raising and capacity-building

Member states are encouraged to:

- ▶ introduce a **right to information on algorithmic mediation** in the context of discrimination complaints or claims.

- ▶ encourage the **rolling out of digital literacy programmes** to raise awareness among citizens of their digital rights relating to equality, including gender equality and non-discrimination.
- ▶ strengthen legal requirements on **democratic participation in standard-setting** given the prominent role that AI standardisation plays in relation to equality, including gender equality and non-discrimination.
- ▶ invest in **capacity-building** including interdisciplinary research on non-discriminatory algorithms and into strategies to protect equality in the use of algorithmic systems.

Introduction: The context

Artificial intelligence (AI) is everywhere. Often acclaimed for its ability to reduce friction and simplify previously manual and time-consuming processes, AI research continues to hurtle down the scientific highway, crossing frontiers and changing the way people live their lives.

In healthcare, the automation of medical diagnosis could make complex services like breast cancer screening and MRI scans function as a walk-in service. This would enable dangerous diseases to be diagnosed in greater volumes and at a much earlier stage. Smart cities can support better management of traffic and allocation of resources, and large-scale data analysis can optimise resources for our environment. AI is also increasingly relied upon as an information and decision-making tool in the world of government and public policy, from housing and healthcare to education and criminal justice. More recently, there has been a lot of discussion in the media around generative AI (a specific type of AI that is focused on generating new content, such as text, images, music etc. using deep learning algorithms like GAN, Transformers and others) due to the accessibility and use of ChatGPT3, an evolution of generative AI which resembles human like conversations and can be used for generating computer code, college-level essays, poems, etc.

Over recent years, the potential to greatly benefit people has been somewhat eclipsed by the growing awareness of a downside: the potential for the *softwarisation*² of existing discrimination and inequality. For example, in what the Dutch have dubbed the “toeslagenaffaire”, or the childcare benefits scandal, thousands of people have suffered the consequence of a biased self-learning algorithm that created risk profiles in an effort to spot childcare benefits fraud. The victims of this case of algorithmic profiling experienced distress and increased poverty, even leading to a case of attempted suicide.³ A parliamentary report into the childcare benefits scandal found several grave shortcomings, including institutional biases and authorities hiding information or misleading the parliament about the facts.⁴

-
2. The “softwarisation” of bias means that existing inequalities end up coded in and perpetuated in obscure and IP-protected machines, see page 10 for further explanation.
 3. Melissa Heikkilä, Dutch scandal serves as a warning for Europe over risks of using algorithms, Politico, 29 March 2022, available at: <https://www.politico.eu/article/dutch-scandal-serves-as-a-warning-for-europe-over-risks-of-using-algorithms/> (last accessed: 30 August 2022)
 4. See Tweede Kamer der Staten-Generaal, Parlementaire ondervraging kinderopvangtoeslag (2020) available at: <https://zoek.officielebekendmakingen.nl/kst-35510-1.pdf>.

In 2018, Reuters reported that Amazon tried to use AI to build a CV-screening tool by using CVs that the company had collected over the previous decade.⁵ As these CVs came mostly from men, and as the consequences of that fact were not seriously thought through, the new system discriminated against women and had to be discarded. In 2019, the Apple-branded credit card came under intense scrutiny because women were receiving less credit than their male spouses who had the same income and credit score.⁶

These cases are neither fringe nor extreme scenarios. Algorithmic systems are too often built and sustained by historic data and models that reproduce stereotypes and false assumptions about gender, race, sexual orientation, ability, class, age, religion or belief, geography, and other socio-cultural and demographic factors. **The bottom line is that without dedicated effort, the use of algorithmic technologies perpetuates and amplifies societal inequalities and harmful stereotypes.**

Awareness of the risks of algorithmic discrimination has crystallised around discussions on ‘bias’, which has now become a prominent public issue. A 2022 survey showed that over 36% of companies “experience[e] challenges or direct business impact due to an occurrence of AI bias in their algorithms, such as [...] [l]ost revenue, [l]ost customers, [l]ost employees, [i]ncurred legal fees due to a lawsuit or legal action [and] [d]amaged brand reputation/media backlash”.⁷ Legislators and regulators around the world are also grappling with these risks and with the pitfalls of existing legislation to address them. Questionnaire answered by the members and observers of the GEC and the CDADI for the purpose of the present Study show broad awareness of the legal issues related to algorithmic bias.⁸ In almost all State Parties, policy or legislative initiatives are either ongoing or public consultations are taking place for this purpose.

The Council of Europe has undertaken work in this area. The Ad Hoc Committee on Artificial Intelligence (CAHAI) was mandated in 2019-2021 to

-
5. Amazon scraps secret AI recruiting tool that showed bias against women, Reuters, 11 October 2018, available at: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G> (last accessed: 25 July 2022).
 6. Alisha Haridasani Gupta, “Are Algorithms Sexist?” *The New York Times* (15 November 2019) available at: <https://www.nytimes.com/2019/11/15/us/apple-card-goldman-sachs.html> (last accessed: 25 July 2022).
 7. See DataRobot, “DataRobot’s State of AI Bias Report Reveals 81% of Technology Leaders Want Government Regulation of AI Bias” (2022), available at: <https://www.datarobot.com/newsroom/press/datarobots-state-of-ai-bias-report-reveals-81-of-technology-leaders-want-government-regulation-of-ai-bias/>.
 8. See section II.

consult with stakeholders and to examine the feasibility and potential elements of a legal framework for the development, design and application of artificial intelligence, based on the Council of Europe standards on human rights, democracy and the rule of law. The Committee published a “Feasibility Study on legal framework on AI design, development and application based on Council of Europe standards” in 2020 as well as “Possible elements of a legal framework on artificial intelligence, based on the Council of Europe’s standards on human rights, democracy and the rule of law”. Following these developments, a new Committee on Artificial Intelligence (CAI) was set up in 2022 and mandated to draft a Framework Convention “on the development, design, and application of artificial intelligence systems based on the Council of Europe’s standards on human rights, democracy and the rule of law, and conducive to innovation”⁹ A legally binding instrument of the Council of Europe has the valuable potential to foster a **human-rights-based approach** to the use of AI and algorithmic technologies in and beyond the international community of State Parties to the ECHR. In addition, the Committee of Ministers has instructed the Gender Equality Commission and the Steering Committee on Anti-Discrimination, Diversity and Inclusion to contribute to the work on a possible legal framework for artificial intelligence systems, to develop a study on the impact of artificial intelligence systems, their potential for promoting equality, including gender equality, and the risks they may cause in relation to non-discrimination, and subject to the results of the study, develop in close co-operation with CAI a possible specific legal instrument.

The aim of this study is threefold. First, it explains how bias in AI and algorithmic technologies arises and may lead to discrimination. It highlights how bias is not just related to data but to the wider human and social underpinnings of these technological artefacts. Second, the Study reviews how policy makers, legislators and companies are dealing with the discriminatory risks of algorithmic technologies and assesses which existing legal instruments could be used for this purpose in the future. It also identifies the shortcomings of existing legal tools and proposes regulatory adaptations to promote equality and prevent discrimination from arising in the development and deployment of algorithmic systems. Third, the Study explores the socio-political conditions necessary for algorithmic technologies to be used to promote equality. It sets out possibilities to leverage these technologies for equality through the legal routes of positive action and positive obligations. Finally, the Study recommends several avenues for ensuring that the use of

9. See Terms of reference of the Committee on Artificial Intelligence CM(2021)131 available at: <https://rm.coe.int/cai-terms-of-reference/1680a7b90b>.

algorithmic technologies does not automate existing inequalities but contributes to a better and more equitable society. All in all, this study aims to support the work of a future Expert Committee under the GEC and CDADI to draft a possible specific sectoral legal instrument on the impact of artificial intelligence systems on equality, including gender equality, and non-discrimination in 2024 and 2025.

In terms of scope, the study focuses mostly on Europe and charts the opportunities and problems that the deployment of algorithmic technologies in society poses in relation to equality and discrimination. It explores the responses that have been given and are being discussed in several countries that are members of the Council of Europe or have observer status to the GEC or the CDADI. The Study builds on Borgesius' study on "Discrimination, Artificial Intelligence and Algorithmic Decision-Making" commissioned by the Council of Europe in 2018 as well as on the fast-developing interdisciplinary body of research on algorithmic discrimination and AI bias.¹⁰ The Study addresses issues of algorithmic discrimination across all grounds protected under Article 14 of the European Convention on Human Rights (ECHR) but with a particular focus on the three groups of protected grounds that are gender and sex, gender identity and sex characteristics; and race, ethnic and national origin, colour, citizenship, religion, language. The study reviews the harmful consequences of AI bias in a wide range of public and private sectors, but with an emphasis on employment and education. Finally, the Study focuses on the legal context and instruments of the Council of Europe, but aligns with and complements the risk-based approach adopted by the European Union in its proposed EU AI Act.

10. See Frederik Borgesius, *Discrimination, Artificial Intelligence and Algorithmic Decision-Making* (2018) Council of Europe available at: <https://rm.coe.int/discrimination-artificial-intelligence-and-algorithmic-decision-making/1680925d73>. See also Fundamental Rights Agency, *Bias in Algorithms – Artificial Intelligence and Discrimination* (2022) available at: https://fra.europa.eu/sites/default/files/fra_uploads/fra-2022-bias-in-algorithms_en.pdf and Carsten Orwat, 'Diskriminierungsrisiken durch Verwendung von Algorithmen' (Antidiskriminierungsstelle des Bundes, 2019).

Section 1

Unpacking ‘machine bias’: How can algorithmic technologies lead to discrimination?

A note on terminology: For the sake of clarity, the term “**user**” of algorithms refers to companies, public bodies or any other stakeholders who deploys an algorithm to support or automate a decision-making process. By contrast, “**end users**” are those subjected to algorithmic or algorithmically supported decisions, for instance customers, job candidates, tax payers, etc. “**Providers**” of algorithmic and AI systems are those who design and commercialise such systems without implementing them in real-life conditions. Sometimes, when algorithmic or AI systems are developed in-house, the provider and the user are the same entity.

1) What is AI?

For the purpose of this analysis, we use the **broad definition of AI** put forward by the ad hoc committee on artificial intelligence (CAHAI) of the Council of Europe, which describes AI “as a ‘blanket term’ for various computer applications based on different techniques, which exhibit capabilities commonly and currently associated with human intelligence”.¹¹ The CAHAI acknowledges that “[t]hese techniques can consist of formal models (or symbolic systems) as well as data-driven models (learning-based systems) typically relying on statistical approaches, including for instance supervised learning, unsupervised learning and reinforcement learning” and that “AI systems act in the physical or digital dimension by recording their environment through data acquisition, analysing certain structured or unstructured data, reasoning on the knowledge or processing information derived from the data, and on that basis decide on the best course of action to reach a certain goal”.¹²

11. Ad hoc committee on artificial intelligence, Feasibility Study CAHAI(2020)23 (Council of Europe, 2020), [8].

12. Ibid.

A further aspect of the definition is that “[these systems] can be designed to **adapt their behaviour over time based on new data** and enhance their performance towards a certain goal”¹³

The background to this broad definition of AI is that, **to date, there is no single definition of AI accepted by the scientific community**. For example, the proposed EU AI regulations define AI as “software that is developed with one or more [...given...] techniques and approaches and can, for a given set of human-defined objectives, generate outputs such as content, predictions, recommendations, or decisions influencing the environments they interact with”¹⁴

According to the EU definition, the techniques and approaches leading to software being identified as an AI system include:

- ▶ “Machine learning (including supervised, unsupervised and reinforcement learning, using a wide variety of methods including deep learning);
- ▶ Logic- and knowledge-based approaches (including knowledge representation, inductive (logic) programming, knowledge bases, inference/deductive engines, (symbolic) reasoning and expert systems);
- ▶ Statistical approaches, Bayesian estimation, search and optimization methods”¹⁵

This variety of techniques falling under the definition of AI include software powering, for example, search engines, image and speech recognition systems, machine translation websites, virtual assistants, spam filters, programmes supporting medical diagnosis, as well as machines such as self-driving cars, robots, and a myriad of objects falling under the vast category of the Internet of Things.¹⁶ In this Study, we find it important to underline that **the regulatory subject is not AI taken in isolation but rather the broader socio-technical apparatus** constituted by the interaction of social elements with algorithmic technologies.

13. Ibid.

14. EU AI Act, Art. 3(1).

15. See Annex 1 of the EU AI Act: “Artificial intelligence techniques and approaches referred to in Article 3, point 1”.

16. European Parliament, “What is artificial intelligence and how is it used?” (2021) available at: <https://www.europarl.europa.eu/news/en/headlines/society/20200827STO85804/what-is-artificial-intelligence-and-how-is-it-used>.

2) What is algorithmic bias?

Algorithms are able to process a far greater range of inputs and variables to make decisions, and can do so with speed and, in many fields, reliability that far exceed human capabilities. From the ads we are served, to the products we are offered, and to the results we are presented with after searching online, algorithms play an ever-greater part in making these decisions.

However, because algorithms simply present the results of calculations **defined by humans** using big data collected from humans, machines, or a combination of the two (at some point during the process), they reflect and process the human biases that are incorporated when the algorithm is programmed, when it processes data and when humans interact with it.¹⁷

In a nutshell, “[algorithmic] **[b]ias happens when seemingly innocuous programming takes on the prejudices either of its creators or the data it is fed.**”¹⁸ As a consequence, women for example (especially from minority groups) may be denied loans and credit, and speech recognition programs may misidentify words spoken by black people at much greater rates than for white people.¹⁹

As Sofiya Noble’s concept of “algorithmic oppression” clarifies, bias is not a “glitch” in otherwise unbiased systems but is instead **systemic and inherent in the functioning of information systems** powering search engines and other web applications.²⁰

Contrary to a widespread narrative, datasets are not the only relays of bias in learning algorithms. Bias has different sources throughout the lifecycle of algorithmic applications, from their inception to their deployment and use. **The complexity of bias emergence and impact is the reason why close attention must be paid to the entire lifecycle of AI and algorithmic systems.**²¹ Several taxonomies listing the sources of bias and its channeling into

17. See e.g., Cathy O’Neil, *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy* (Broadway Books 2017).

18. Garcia, Megan. “Racist in the Machine: The Disturbing Implications of Algorithmic Bias.” *World Policy Journal* 33 (2016): 111 - 117.

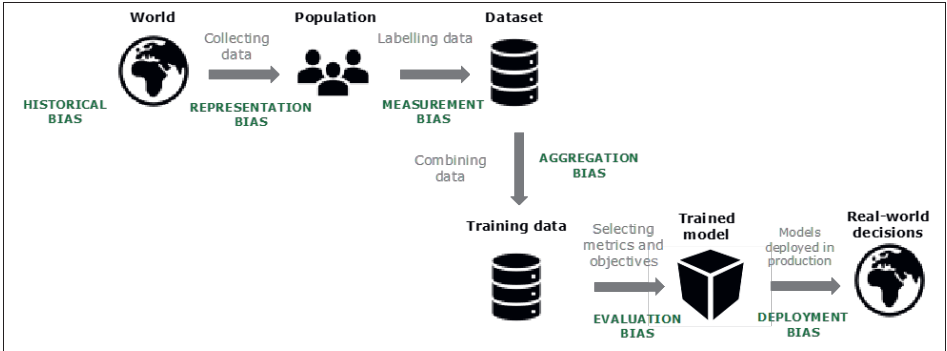
19. Allison Koenig, et al., PNAS, March 23, 2020

20. See Safiya Noble, *Algorithms of Oppression: How Search Engines Reinforce Racism* (New York: New York University Press, 2018) and Vanessa Ceia, Benji Nothwehr, and Liz Wagner, *Gender and Technology: A rights-based and intersectional analysis of key trends* (Oxfam Research Background, 2021), 40.

21. Ivana Bartoletti, *The Complex Issue of Algorithmic Fairness*, The Yuan, September 2021, available at: <https://www.the-yuan.com/129/The-Complex-Issue-of-Fairness-in-AI-Part-I.html> (last accessed: 28 July 2022)

AI systems and outputs have been developed by researchers. For example, the diagram below by Suresh and Guttag shows the different entry points for bias, and what they entail.

Table and definitions below from: A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle²²



Suresh and Guttag distinguish five sources and types of bias in AI systems. First, what they call “**historical bias**” describes how social hierarchies and institutionalised disadvantage shape social data.²³ Data is therefore not neutral because it reflects the unequal society in which we live. For example, as women have historically earned less than men, they may be given less credit²⁴ or, in the context of advertising, be served adverts for lower paid job posts.²⁵

In turn, “**representation bias**” arises in data collection.²⁶ For example, if an organisation’s marketing team advertises in predominantly white neighbourhoods, the resulting customer base would not be representative of the

22. Harini Suresh and John Guttag. 2021. A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. In Proceedings of EAAMO '21: Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '21). ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3465416.3483305>.
23. See *ibid*.
24. Apple’s ‘sexist’ credit card investigated by US regulator, BBC, 11 November 2019, available at: <https://www.bbc.com/news/business-50365609> (last accessed: 15 June 2022).
25. Samuel Gibbs, Women less likely to be shown ads for high-paid jobs on Google, study shows, The Guardian, 8 July 2015, available at: <https://www.theguardian.com/technology/2015/jul/08/women-less-likely-ads-high-paid-jobs-google-study> (last accessed: 15 June 2022).
26. See Harini Suresh and John Guttag. 2021. A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. In Proceedings of EAAMO '21: Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '21). ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3465416.3483305>.

wider population. That dataset would generate bias if used for example to train an algorithm later used to cater to broader population groups.²⁷

The researchers also shed light on “**measurement bias**”, which “occurs when choosing, collecting, or computing features and labels to use in a prediction problem”.²⁸ Many features and labels are non-problematic, such as the labelling of an image as a cat or a dog, but problems may emerge when some factors are used as a proxy. For example, postcode could be a proxy for race or sexual orientation, occupation could be a proxy for gender and first names are often used as proxies for age.²⁹ Alternatively, if proxies overly simplify the feature to be measured or the proxy reflects variations in the quality of measurements across groups, measurement bias could arise.³⁰

“**Aggregation bias**” relates to how data is combined. It occurs when data groups are inappropriately combined, resulting in a model that does not perform well for any group or only performs well for the majority group.³¹ The researchers mention the example of local meanings ascribed by specific communities to emoji, hashtags and sentences on social media, which differ from the meanings in the broader social media user population.³² This could lead for instance to content moderation applying inadequate semantic filters modelled on majority groups to minority groups, with silencing effects that could unfairly restrict minority groups’ ability to communicate via social media.

The researchers also identify “**evaluation bias**”, which occurs when evaluating a model, if the benchmark data (used to compare the model to other models that perform similar tasks) does not represent the population that the model will serve.³³ For example, the Gender Shades paper discovered

27. See further the examples on p. 15 in relation to men being used as the baseline for health-care research, in Criado Perez C, *Invisible women: Exposing data bias in a world designed for men* (Random House 2019).

28. *Ibid.*

29. See various tools that are designed to predict age from data about names: <https://cebus.net/de/age.php>: <https://agify.io/> or <https://github.com/JasonKessler/agefromname>.

30. *Ibid.*

31. *Ibid.*

32. *Ibid.*, citing a study by Desmond U. Patton, William R. Frey, Kyle A. McGregor, Fei-Tzin Lee, Kathleen McKeown, and Emanuel Moss. 2020. Contextual Analysis of Social Media: The Promise and Challenge of Eliciting Context in Social Media Posts with Natural Language Processing. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (New York, NY, USA) (AIES '20). Association for Computing Machinery, New York, NY, USA, 337–342. <https://doi.org/10.1145/3375627.3375841>.

33. See Harini Suresh and John Guttag. 2021. A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. In Proceedings of EAAMO '21: Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '21). ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3465416.3483305>.

that two widely used facial analysis benchmark datasets (IJB-A and Adience) were primarily composed of lighter-skinned subjects (79.6% and 86.2%, respectively).³⁴

Finally, “**deployment bias**” relates to the real-world use of models, in particular if a model developed to solve a problem is used for another task.³⁵ This could happen for example due to a change in marketing strategy. In addition, a model is often a part of a complex socio-technical system where humans and machines interact. In a ‘live’ environment, additional biases may therefore be introduced when humans interpret algorithmic outputs to be used as inputs further down the algorithmically supported decision-making line.³⁶

So-called **automation and confirmation biases** can also strengthen these biases. Automation bias takes place when humans place greater trust in machines and technological artefacts than in their own or other humans’ potentially contradictory judgment, and therefore tend to validate algorithmic outputs without questioning them. In the context of predictive machines for example, such bias can lead to biased risk assessments not being challenged by so-called humans-in-the-loop and therefore to rubberstamping behaviours. Confirmation bias happens when pre-existing beliefs influence the processing of new information, leading in particular to new information being better retained when consistent with such beliefs or being interpreted in consistency with such beliefs. In the AI context, this could lead to gender stereotypes acting as a reinforcing prism by human decision-makers when interpreting biased algorithmic outputs. In an experiment, Green and Chen also shows that human interpreters of automated risk assessments provided by an algorithm yield “**disparate interactions**”, that is interpretations of similar algorithmic risk assessments are more lenient towards white than black defendants.³⁷

Other taxonomies of bias have been proposed. For example, Barocas and Selbst identify key moments and situations where bias is channeled into AI systems: the **definition of “target variables”** (the feature to be measured or

34. Buolamwini J and Gebru T, *Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification* (Proceedings of Machine Learning Research 2018).

35. See *ibid.*

36. Harini Suresh and John Guttag. 2021. A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. In Proceedings of EAAMO '21: Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '21). ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3465416.3483305>.

37. See Green B and Chen Y, ‘Disparate interactions: An algorithm-in-the-loop analysis of fairness in risk assessments’ (2019) Proceedings of the Conference on Fairness, Accountability, and Transparency 90.

predicted by a model, e.g., work performance) and “class labels” (the possible variations in the occurrence of the target variable, for example stellar, very good, good, unsatisfactory); the use of “training data” (with bias occurring during labelling and data collection); “feature selection” (the attributes that are to be considered relevant by a model, for instance yearly income); and the use of “proxies” (when relevant attributes correspond to protected groups, for example yearly income and gender due to the gender pay gap).³⁸

The most recent example of ChatGPT3 provides a clear example of the wide-ranging shades of bias. First, these models are trained on Wikipedia which is a largely male dominated platform. It is worth noting that, for example, English-language Wikipedia contains more than 1.5 million biographies about notable writers, inventors, and academics, but less than 19% of these biographies are about women.³⁹ Questions also arise regarding the diversity of the workforce labelling the data. Finally, ideas circulating about customised models replacing one-size-fits-all ChatGPT to align with our own politics, raise serious issues about human rights and their universalism.⁴⁰

These taxonomies help debunk the myth that bias emerges from data only and show the complex role of socio-technical interactions in the (re) production of discriminatory bias.

3) The discriminatory impact of AI: some concrete examples

This section illustrates how bias can give rise to discrimination across different sectors.

Recruitment: Reuters reported in 2018 that Amazon developed a program relying on machine-learning to identify top candidates in pools of CVs. The program systematically disadvantaged women’s CV because it reflected the gender gap in the workforce recruited over the past ten years. Neutralising

38. Barocas S and Selbst AD, ‘Big Data’s Disparate Impact’ (2016) 104 California law review , 677-693.

39. Tripodi, F. (2021). Ms. Categorized: Gender, notability, and inequality on Wikipedia. *New Media & Society*, 0(0). <https://doi.org/10.1177/14614448211023772>

40. Eric Hal Schwartz, OpenAI Promises Customizable ChatGPT After Bias Complaints, 20 February 2023, available at: <https://voicebot.ai/2023/02/20/openai-promises-customizable-chatgpt-after-bias-complaints/>

words like “women” did not redress the discriminatory outcome as the system was able to infer sex from other data.⁴¹

Researchers based at Utrecht University partnered with a job matching platform to research how the use of gendered language in the search bar yields different results, with discriminatory allocations of information about job opportunities.⁴² This not only results in strengthening stereotypes about male and female typical occupations but also results in allocative and distributive harms.

The online targeted distribution of job adverts powered by optimisation services offered by social media platforms such as Facebook also serves to reinforce gender stereotypes as well as gender segregation within the workplace.⁴³ An experiment conducted by AlgorithmWatch in 2020 showed that when asking Facebook to distribute ads “neutrally” (without targeting a specific audience), an ad for a truck driver position was shown to a public composed of 93% men and 7% women.⁴⁴ Conversely, an advert for a position as educator was distributed to an audience composed of 96% women and 4% men.⁴⁵

AI-powered face recognition and emotions analysis systems can also yield racial discrimination or disadvantage job candidates with disabilities.⁴⁶ This

-
41. See Destin J, ‘Amazon scraps secret AI recruiting tool that showed bias against women’ *Reuters* (2018) available at: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G> (last accessed 22 July 2022).
 42. See van Es K, Everts D and Muis I, ‘Gendered language and employment Web sites: How search algorithms can cause allocative harm’ (2021) 26 *First Monday* available at: <https://journals.uic.edu/ojs/index.php/fm/article/view/11717/10200>.
 43. See Ali M and others, ‘Discrimination through optimization: How Facebook’s Ad delivery can lead to biased outcomes’ (2019) 3 *Proceedings of the ACM on Human-Computer Interaction* 1.
 44. 4,864 men, but only 386 women. See Wulf J, *Automated Decision-Making Systems and Discrimination: Understanding causes, recognizing cases, supporting those affected* (AlgorithmWatch 2022), 7 available at: https://algorithmwatch.org/en/wp-content/uploads/2022/07/AutoCheck-Guidebook_ADM_Discrimination_EN-AlgorithmWatch_June_2022_b.pdf and Kayser-Bril N, ‘Automated Discrimination: Facebook uses gross stereotypes to optimize ad delivery’ *AlgorithmWatch* available at: <https://algorithmwatch.org/en/automated-discrimination-facebook-google/> (last accessed 22 July 2022).
 45. *Ibid.* 6,456 women, but only 258 men.
 46. See Buolamwini J and Gebru T, *Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification* (*Proceedings of Machine Learning Research* 2018); Hannah Devlin, “AI systems claiming to ‘read’ emotions pose discrimination risks” (16 February 2020) *The Guardian* available at: <https://www.theguardian.com/technology/2020/feb/16/ai-systems-claiming-to-read-emotions-pose-discrimination-risks> (last accessed 22 July 2022).

is because of lower performance rates of such devices on darker skin tones, especially for women.⁴⁷ In addition, emotions analysis software trained on neurotypical subjects might not be able to perform correctly on neurodiverse subjects. As AI-powered emotions analysis is increasingly used in the recruitment sector, for instance to analyse video recordings of job candidates' presentations, this could pose accessibility and inclusion issues.

Access to goods and services, banking and insurance: In Finland the National Non-Discrimination and Equality Tribunal found direct multiple discrimination in a case where the applicant was denied a loan online. After investigating the case, the Equality Body (the Non-Discrimination Ombudsman) found that the company had used statistical models to assess credit worthiness that relied on an applicant's age, gender, language and place of residence while not taking into account an applicant's actual credit history. In that case, the applicant being male, Finnish speaker and from a rural area were treated as factors of disadvantage in the assessment performed by the financial institution.⁴⁸

A similar story was reported in Germany, where a female customer was refused credit while purchasing goods online. When investigating the reasons for the rejection with the credit institution, the customer learned that a combination of her age and gender seemed to have motivated the automated rejection, based on harmful intersectional stereotypes that women around 40 are often divorced and have therefore less economic power.⁴⁹

In the insurance sector, a study conducted by the Universities of Padua, Udine, and Carnegie Mellon showed that factors such as birthplace and citizenship influence the price of car insurance policies paid by customers.⁵⁰ In a case study, they showed that indicating Ghana as an applicant's birthplace could lead to a price increase of 1000 EUR compared to an applicant indicating Italy as their birthplace.

47. See Buolamwini J and Gebru T, Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification (Proceedings of Machine Learning Research 2018).

48. See Lorenz Matzat and Minna Ruckenstein, "Finnish Credit Score Ruling raises Questions about Discrimination and how to avoid it" (21 Novembre 2018) *AlgorithmWatch* available at: <https://algorithmwatch.org/en/finnish-credit-score-ruling-raises-questions-about-discrimination-and-how-to-avoid-it/> (last accessed 22 July 2022); Rainer Hiltunen, "Multiple discrimination in assessing creditworthiness" (1 August 2018) European network of legal experts in gender equality and non-discrimination available at: <https://www.equalitylaw.eu/downloads/4658-finland-multiple-discrimination-in-assessing-creditworthiness-pdf-120-kb> (last accessed 22 July 2022).

49. See Wulf J, Automated Decision-Making Systems and Discrimination: Understanding causes, recognizing cases, supporting those affected (AlgorithmWatch 2022), 6-7

50. The study was reported by AlgorithmWatch, see *ibid*.

Another study by AlgorithmWatch showed that digital discrimination extends far beyond AI.⁵¹ Simple online forms can cause discrimination on grounds of race, ethnic origin or nationality, for example if they only allow registering patronyms containing three or more letters. Applicants with shorter names will be denied registration or unable to open an account, which is often a precondition for purchasing goods and services online.

Risk assessment in the area of security, crime prevention, policing and the justice system: In Spain the VioGén software has been used to assess risks of gender-based violence and femicide by intimate partners. Despite an overall favourable assessment, criticisms point to several cases of false negatives where low risk scores led to insufficient prevention means being deployed, with tragic consequences.⁵²

The Netherlands have deployed several predictive systems for crime prevention purposes, which have been harshly criticised for creating discrimination based on race, ethnicity and nationality. For instance, a 2020 investigation by Amnesty International revealed that the “Sensing Project”, that aimed to prevent shoplifting and pickpocketing locally, resulted in discriminatory ethnic profiling of individuals of Eastern European origin, and in particular members of the Roma community.⁵³ When surveilling car traffic in and around the area of deployment, the system used the Eastern European origin of passengers as a predictive risk factor for crime. Other crime anticipation systems, for instance in Amsterdam, have been reported to use factors such as “number of one parent households”, “number of social benefits recipients” and “number of non-Western immigrants” to identify crime “hot spots” throughout the country⁵⁴

At airports, security screening and border control technologies using automated gender recognition systems have been shown to discriminate against transgender, intersex, non-binary and gender non-conforming persons

-
51. Lulamae, Josephine, “Fixing Online Forms Shouldn’t Wait Until Retirement”, AlgorithmWatch (13 January 2022) available at: <https://algorithmwatch.org/en/unding-online-forms/> (last accessed 22 July 2022).
 52. Michele Catanzaro, “In Spain, the VioGén algorithm attempts to forecast gender violence”, AlgorithmWatch (27 April 2020) available at: <https://algorithmwatch.org/en/viogen-algorithm-gender-violence/> (last accessed 22 July 2022).
 53. Amnesty International “We Sense Trouble: Automated Discrimination and Mass Surveillance in Predictive Policing in the Netherlands” (2020), 5 available at: https://www.amnesty.nl/content/uploads/2020/09/Report-Predictive-Policing-RM-7.0-FINAL-TEXT_CK-2.pdf (last accessed 22 July 2022).
 54. <https://www.vice.com/en/article/5dpmdd/the-netherlands-is-becoming-a-predictive-policing-hot-spot>

because they rely on a binary gender classification system that does not capture the real complexity of gender identity and gender expression.⁵⁵

Facial recognition is increasingly deployed for crime detection and prevention. For example, law enforcement agencies may use face recognition to compare suspects' photos to mugshots and driver's license images. While "[f]ace recognition algorithms boast high classification accuracy (over 90%)", these outcomes are not universal.⁵⁶ In 2018, the Gender Shades project revealed discrepancies in the classification accuracy of face recognition technologies for different skin tones and sexes. These algorithms consistently demonstrated the poorest accuracy for darker-skinned females and the highest for lighter-skinned males.⁵⁷ In a criminal justice setting, face recognition technologies that are inherently biased in their accuracy can potentially misidentify suspects and even lead to the incarceration of innocent people of colour as has happened in the US.⁵⁸ It is therefore concerning that, even if accurate, "face recognition empowers [...] law enforcement system[s] with a long history of racist and anti-activist surveillance and can widen pre-existing inequalities".⁵⁹

Access to public and administrative services: the use of face recognition technologies within or in association with public services can lead to excluding or denying end users public services. For instance, a photo booth at the State Office of Transportation in Hamburg, Germany, failed to recognise an applicant's face for the purpose of taking a biometric picture, which was needed for her administrative application. Even though the public office denied that the failure stemmed from the facial recognition software used,

55. See JD Shadel, "#TravelingWhileTrans: The trauma of returning to 'normal'" (The Washington Post, 2021) available at: <https://www.washingtonpost.com/travel/2021/06/16/trans-travel-tsa-lgbtq/> and Quinan, C. L., and Mina Hunt. "Biometric Bordering and Automatic Gender Recognition: Challenging Binary Gender Norms in Everyday Biometric Technologies." *Communication, Culture and Critique* 15.2 (2022): 211-226.

56. Alex Najibi, *Racial Discrimination in Face Recognition Technology*, Harvard University, October 2020, available at: <https://sitn.hms.harvard.edu/flash/2020/racial-discrimination-in-face-recognition-technology/#:~:text=Face%20recognition%20algorithms%20boast%20high,and%2018%2D30%20years%20old>

57. Gender Shades Project, available at <http://gendershades.org/overview.html> (last accessed: 31 August 2022)

58. RACE AND WRONGFUL CONVICTIONS IN THE UNITED STATES, available at: https://www.law.umich.edu/special/exoneration/Documents/Race_and_Wrongful_Convictions.pdf (last accessed: 31 August 2022).

59. Alex Najibi, *Racial Discrimination in Face Recognition Technology*, Harvard University, October 2020.

a local employee indicated that failures often take place in relation to applicants' skin colour.⁶⁰

In the Netherlands, the deployment of the SyRi system (System Risk Indication), used to detect social welfare fraud, was shown to cause discrimination on grounds of income and ethnic origin before being put to halt by a court decision in 2020.⁶¹ In 2021, a welfare scandal forced the Dutch government to resign after more than 20.000 parents were flagged by an AI system as fraudsters in relation to childcare allowance and subjected to investigation by the Dutch tax authorities.⁶² The AI system treated double nationality as a high risk factor and this resulted in a disproportionate number of investigations and court proceedings being launched against families with an immigration background, whose child care benefits were suspended and some of whom were requested to reimburse the benefits received.⁶³ The case also shows how the lack of accountability and transparency around the use of these systems can lead to depriving the subjects of AI decision-making from an explanation or the opportunity to appeal against the decisions.

Education: Facial recognition software have been known to be biased and lead to intersectional discrimination on grounds of race and gender.⁶⁴ When used in proctoring software in educational settings, that can negatively affect the conditions in which racialised students take exams and even their ability to do so. For example, proctoring software used by several universities in

-
60. See Wulf J, Automated Decision-Making Systems and Discrimination: Understanding causes, recognizing cases, supporting those affected (AlgorithmWatch 2022), p8. This hypothesis is corroborated by studies pointing at intersectional discrimination on grounds of gender and skin colour in facial recognition software, e.g., Buolamwini J and Gebru T, Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification (Proceedings of Machine Learning Research 2018).
 61. Koen Vervloesem, "How Dutch activists got an invasive fraud detection algorithm banned", AlgorithmWatch (6 April 2020) available at: <https://algorithmwatch.org/en/syri-netherlands-algorithm/> (last accessed 22 July 2022).
 62. Nadia Benaissa, "Het systeem doet precies wat het wordt opgedragen" (29 January 2021) Bits of Freedom available at: <https://www.bitsoffreedom.nl/2021/01/29/het-systeem-doet-precies-wat-het-wordt-opgedragen/>.
 63. Jon Henley, "Dutch government faces collapse over child benefits scandal" (14 January 2021) The Guardian available at: <https://www.theguardian.com/world/2021/jan/14/dutch-government-faces-collapse-over-child-benefits-scandal> and Björn ten Seldam & Alex Brenninkmeijer, "The Dutch benefits scandal: a cautionary tale for algorithmic enforcement" (30 April 2021) EU Law Enforcement available at: <https://eulawenforcement.com/?p=7941>.
 64. Buolamwini J and Gebru T, Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification (Proceedings of Machine Learning Research 2018).

the Netherlands had trouble recognising dark-skinned students.⁶⁵ After the University did not take her complaint seriously, a student supported by the Racism and Technology Center submitted a formal complaint to the Institute of Human Rights, the national non-discrimination authority in the country.⁶⁶ Proctoring software can also impact students with disabilities negatively, for instance by generating anxiety, not allowing a carer or not letting the students take breaks away from the computer.⁶⁷ For low-income families who share rooms due to a lack of space at home, the use of a proctoring software can create disadvantage by signalling “aberrant behaviour” if family members are identified passing behind the screen.⁶⁸

Healthcare: Criado Perez has exposed how healthcare research and industry rely on male models to assess the risks and efficacy of drugs, thus yielding less and lower quality health data for women and gender diverse persons. Such gender data gap in the healthcare sector, leads to less reliable predictive systems when it comes to diagnosing female and gender diverse patients.⁶⁹ Research shows that the data gap in health also affects other minority groups.⁷⁰

A US study by Obermeyer et al. shows how a system used to predict health-related risks in order to allocate resources systematically disadvantaged patients with ethnic minority backgrounds. This is because the system used

-
65. Racism and Technology Centre, “Student stapt naar College voor de Rechten van de Mens vanwege gebruik racistische software door de VU” (15 July 2022) available at: <https://racismandtechnology.center/2022/07/15/student-stapt-naar-college-voor-de-rechten-van-de-mens-vanwege-gebruik-racistische-software-door-de-vu/#more-1691> (last accessed 28 July 2022).
 66. Fleur Damen, “De antispieksoftware herkende haar niet als mens omdat ze zwart is maar bij de vu vond ze geen gehoor” De Volkskrant (15 July 2022) available at: <https://www.volkskrant.nl/nieuws-achtergrond/de-antispieksoftware-herkende-haar-niet-als-mens-omdat-ze-zwart-is-maar-bij-de-vu-vond-ze-geen-gehoor~b6810279/> (last accessed 27 July 2022). See the complaint at: <https://racismandtechnology.center/2022/07/15/student-stapt-naar-college-voor-de-rechten-van-de-mens-vanwege-gebruik-racistische-software-door-de-vu/#more-1691>
 67. Lydia X. Z. Brown, “How Automated Test Proctoring Software Discriminates Against Disabled Students” (16 November 2020) Centre for Democracy and Technology available at <https://cdt.org/insights/how-automated-test-proctoring-software-discriminates-against-disabled-students/> (last accessed 28 July 2022).
 68. Ibid.
 69. See Criado Perez C, *Invisible women: Exposing data bias in a world designed for men* (Random House 2019).
 70. Ibid.

data about groups' previous access to healthcare, which embedded existing structural discrimination.⁷¹

A Study published by the World Health Organisation in 2022 shows that algorithmic systems used in the healthcare sector are trained on the data of predominantly younger populations, which is not representative of older subjects.⁷² This decreases the quality of predictions for older populations and could lead to disproportionately lower performance of these systems, including with incorrect diagnosis.

Media and search engines: Research shows that representations of women in images returned by search engines online are biased and reflect sexist, racist and intersectionally discriminatory stereotypes. For instance, Noble shows in an experiment with the Google search engine how images of black girls and black women are sexualised.⁷³ Other groups of minority women are also subjected to sexualised stereotyping in search engines results, for example in searches related to the word "lesbian".⁷⁴ Even though search engines have tried to correct these biases, a recent study surveying major search engines shows "representation bias" as well as "face-ism bias" in the way in which women are represented, meaning that "[w]omen are less likely to be represented in gender-neutral media content [...] and their face-to-body ratio in images is often lower" than for men.⁷⁵ Technical debiasing solutions might treat some of the symptoms of the problem, for instance re-balancing the amount of female pictures in an image search for "CEOs", but not its roots,

71. See Obermeyer Z and others, 'Dissecting racial bias in an algorithm used to manage the health of populations' (2019) 366 *Science* 447.

72. J Stypinska, 'AI ageism: a critical roadmap for studying age discrimination and exclusion in digitalized societies' (2022) *AI & Soc* available at: <https://doi.org/10.1007/s00146-022-01553-5> and WHO, 'Ageism in artificial intelligence for health' (2022) available at: <https://www.who.int/publications/i/item/9789240040793>.

73. See e.g., Safiya Noble, *Algorithms of oppression : how search engines reinforce racism* (New York University Press 2018).

74. The prevalence of misogynistic conceptions concerning non-heterosexual women results in the fact that in many languages there is an association between the word "lesbian" and pornographic contents. Search engines' algorithms replicate such association and influence the type of results obtained by searching such words. For example, in 2019, Google had to change its algorithms to avoid that search results associated with the word "lesbienne" yielded results only linked with pornographic content, while other words associated with the LGBTI community did not bring the same results. See Marie Turcan, 'Pourquoi le mot « lesbienne » sur Google ne renvoie-t-il que vers des sites pornographiques ?' (Numerama, 2019) available at: <https://www.numerama.com/politique/478663-pourquoi-le-mot-les-bienne-sur-google-ne-renvoie-t-il-que-vers-des-sites-pornographiques.html>.

75. Ulloa R and others, 'Representativeness and face-ism: Gender bias in image search' (2022) *New Media & Society*.

in this case harmful stereotyping, representational and allocative harms as well as structural inequality that are deeply entrenched in our cultural and material reality. For instance, recent tests seem to show that the AI-powered art tool DALLE2 adds 'diversity prompts' to unspecific queries, for example adding the labels "black" or "female" to a prompt asking the software to generate an image of 'a CEO'.⁷⁶ This approach is analogous to a form of positive action like quotas. It can be criticised for not addressing the roots of such discrimination, namely the lack of diversity in training sets, but if used at a large scale, such fixes have at least the merit to disseminate more diverse representations that, in the long run, can contribute to mitigating harmful stereotypes.

Online gender-based violence, hate speech, harassment: Digital discrimination also takes the form of gender-based violence, for instance when deep-fake videos are used to harass women in the context of so-called "revenge porn" cases. Unconsented dissemination of sexual content, often in the form of images, has also been recognised as a form of gender-based violence that especially affect women and girls who are young or public figures such as journalists, human rights defenders, or politicians.⁷⁷ In addition, sexist and other forms of online hate speech have been highlighted as contingent on the rising use of social media platforms.⁷⁸ At the same time, content moderation particularly affects minority groups, who are at risk of being silenced⁷⁹ while at the same time subjected to hate campaigns. For example, the stereotypical association of words associated with the lesbian community (e.g., 'lesbian') with pornographic content often results in so-called 'shadow-bans' that limit the reach of social media posts, or in the outright impossibility to use certain words in account names and handles. The silencing effects of content moderation have a severe negative impact on the visibility and

-
76. Matthew Sparkes, "AI art tool DALL-E 2 adds 'black' or 'female' to some image prompts" (22 July 2022) New Scientist available at: <https://www.newscientist.com/article/2329690-ai-art-tool-dall-e-2-adds-black-or-female-to-some-image-prompts/> (last accessed 28 July 2022); see also OpenAI, "Reducing Bias and Improving Safety in DALL-E 2" (18 July 2022) available at: <https://openai.com/blog/reducing-bias-and-improving-safety-in-dall-e-2/> (last accessed 28 July 2022).
 77. See Sara De Vido and Lorena Sosa, Criminalisation of gender-based violence against women in European States, including ICT-facilitated violence (European Network of Legal Experts in gender equality and non-discrimination 2021) available at: <https://www.equalitylaw.eu/downloads/5535-criminalisation-of-gender-based-violence-against-women-in-european-states-including-ict-facilitated-violence-1-97-mb> (last accessed 23 July 2022).
 78. See Bartoletti, Ivana. Chapter 3: Algorithms and the Rise of Populism in An artificial revolution: On power, politics and AI. Black Spot Books, 2020.
 79. See Rachel Griffin, 'The Sanitised Platform' (2022) 13 J Intell Prop Info Tech & Elec Com L 36.

reach of organisations, activities and events aimed at countering hateful and discriminatory narratives targeting such minority communities.⁸⁰

Gender stereotyping across the board: A recent UN report, “I’d blush if I could: closing gender divides in digital skills through education” found that AI digital assistants with female voices can reinforce existing gender biases. This trend toward female voiced virtual assistants “seems to have less to do with sound, tone, syntax, and cadence, than an association with assistance”.⁸¹ Perhaps a female voice is chosen to seduce a user into thinking that AI is pliable and benign. But the ultimate effect is the “normalisation of this new digital servitude in our homes and daily lives through Alexa, Siri and Cortana”.⁸²

4) What makes algorithmic discrimination different?

Discrimination powered by algorithmic technologies presents a set of **distinct challenges**.

First, the deployment of algorithmic systems in decision-making processes entails **large-scale effects on society**. For example, while a bank employee might unconsciously assign a higher mortgage rate to an applicant from a minority group, a software processing thousands of files per day might generalise this bias to any applicant with an African sounding name.

Secondly, human conduct is controlled by social and legal mechanisms that, although far from perfect, are meant to correct misbehaviours in the short and long term. By contrast, **the deployment of algorithmic technologies often jeopardises accountability for, transparency in and scrutiny of decision-making processes**. For example, whereas human decisions can in principle be appealed, the lack of information about AI deployment, the opacity of the systems used and the unwillingness of providers to open up such algorithmically supported decision-making processes to public scrutiny

80. For example, the Eurocentralasian Lesbian Community (EL*C), an organisation created in 2017 to advocate for the rights of LGBTI women, reports that it was unable to use the word “lesbian” in its username on Facebook whereas other words associated with the LGBTI community (such as “gay” or “queer”) could be used. See EL*C, ‘Lesbophobia: An intersectional form of violence’ (2021): <https://europeanlesbianconference.org/wp-content/uploads/2021/10/Lesbophobia-3.pdf>

81. Unesco, I’d blush if I could: closing gender divides in digital skills through education, 100 available at: <https://unesdoc.unesco.org/ark:/48223/pf0000367416.page=1>

82. Ivana Bartoletti, An Artificial Revolution: on Power, Politics and AI (Indigo Press).

make accountability difficult to achieve in AI systems.⁸³ These hurdles also make the collection of evidence of algorithmic discrimination difficult.

Third, **the sources of algorithmic discrimination are difficult to identify.** Due to the complexity and dynamic features of these socio-technical systems, bias can affect any stage of the algorithmic pipeline. In addition, **algorithms might be proprietary, complex and difficult to understand.** Sometimes, they are effectively a sealed box, containing proceedings that may be unexplainable to a human researcher. This “softwarisation” of bias means that existing inequalities end up coded in and perpetuated in obscure and IP-protected machines. This is extremely problematic as bias becomes more difficult to identify and harder to challenge.

To sum up, at least **six challenges** arise with algorithmic and data-driven discrimination.⁸⁴ Machine-supported decisions are made at a much **greater scale** but the interaction between humans and machines make the **sources of discrimination difficult to identify and address.** The ‘**cleaning**’ of **biased data is a technical challenge** and a **context-dependent** exercise, and the existence of proxies for and correlations with protected groups further complicates the task. **Algorithmic determinism** is particularly problematic in relation to discrimination as predictive systems use **correlations arising from historical discrimination** (e.g., the gender pay gap) as **quasi ‘causal’ bases for decision-making**, thereby creating feedback loops. At the same time, AI and algorithmic systems are often **non-transparent**, might not be explainable, and the **attribution of responsibility for discrimination is unclear.**

Because **the source of these biases is not ultimately technological, they cannot be resolved using technology alone.** Instead, addressing algorithmic discrimination and data-driven disadvantage requires a much greater degree of scrutiny and a **positive political decision to actively prevent the reinforcing of structural inequalities engrained in social data.** For example, to avoid “automating” gender stereotypes and the gender pay gap – the fact that women have historically earned less than men – employers need to make a conscious decision to target women when advertising higher paying, typically “masculine” or management jobs online. Simply entrusting their distribution to optimisation algorithms instead is likely to reproduce

83. See Gabriele Spina Ali & Ronald Yu, Artificial Intelligence between Transparency and Secrecy: From the EC Whitepaper to the AIA and Beyond, *European Journal of Law and Technology*, available at: <https://www.ejlt.org/index.php/ejlt/article/download/754/1044/3716> (last accessed: 16 September 2022)

84. See Gerards J and Xenidis R, *Algorithmic discrimination in Europe: Challenges and Opportunities for EU Gender Equality and Non-Discrimination Law* (European network of legal experts in gender equality and non-discrimination / European Commission, 2021).

gender stereotypes and pay inequality.⁸⁵ Understanding algorithmic bias therefore starts with recognising how algorithmic technologies escalate, entrench, and perpetuate existing inequalities where no safeguards are put in place. For these reasons, **addressing algorithmic discrimination requires a multifaceted approach encompassing various disciplines** such as social science, ethics and law, and **regulatory fields** including legislation on non-discrimination, consumer protection, data protection, trade, etc.

5) Addressing algorithmic discrimination: best practices and their limits

To address the discriminatory risks of algorithmic technologies, the industry has taken initiatives ranging from **technical solutions to 'debias' and 'audit' algorithmic systems** to **voluntary codes of conduct, instruments for ethical AI** and other forms of **self-regulation**. This section exhibits some **examples of the good governance practices** adopted and assesses their **limits**.

Companies have been ramping up governance milestones in anticipation of incoming regulation especially as both ex-ante and ex-post governance measures gain popularity and significance. Large tech companies (often themselves hit by controversies around bias) have introduced ethics boards, built AI governance around existing governance structures and/or deployed debiasing techniques to address some of the issues.

For example, Microsoft has developed six AI principles to accelerate this cultural shift and to improve employees' awareness of ethical issues.⁸⁶ These include fairness, reliability and safety, privacy and security, inclusiveness, transparency and accountability. Governance is constituted by three core teams with the purposes of enacting the core principles, management of policy, governance, enablement, and sensitive use functions, and leading the implementation of responsible AI processes in the adoption of systems and tools.

IBM has developed and implemented AI Fairness 360,⁸⁷ an open-source toolkit used to examine, report, and mitigate discrimination and bias in machine

85. See Ali M and others, 'Discrimination through optimization: How Facebook's Ad delivery can lead to biased outcomes' (2019) 3 Proceedings of the ACM on Human-Computer Interaction 1 and Imana B, Korolova A and Heidemann J, *Auditing for discrimination in algorithms delivering job ads* (2021).

86. Microsoft AI Principles, available at: <https://www.microsoft.com/en-us/ai/responsible-ai?activetab=pivot1%3aprimarary6> (last accessed: 4 October 2022)

87. IBM, introducing AI Fairness 360, available at: <https://www.ibm.com/blogs/research/2018/09/ai-fairness-360/> (last accessed: 4 October 2022).

learning models. The main objectives of this toolkit are to help facilitate the transition of fairness research algorithms for use in an industrial setting and to provide a common framework for fairness researchers to share and evaluate algorithms.

Amazon has integrated new tools to assist in detecting discrimination in AI and ML technologies. As part of the cloud computing offering Amazon Web Services, a new test has been implemented alongside a wider suite of materials to customers seeking to develop fair, non-biased AI on the platform. The test was developed by Wachter, Mittelstadt and Russell from the Oxford Internet Institute of the University of Oxford and it is called 'the Conditional Demographic Disparity (CDD)'; a new test for "ensuring fairness in algorithmic modelling and data driven decisions".⁸⁸

The developers of the image generation AI 'DALLE-2' have implemented a bias mitigation technique after evidence of representational harm in image outputs mounted. While generic prompts such as 'CEO' and 'builders' mostly generated images of men, prompts such as 'flight attendant' and 'nurse' generated images representing almost exclusively women.⁸⁹ The developers acknowledge how such stereotypes can be harmful, for instance when harming the dignity of protected groups, erasing them from socially valued situations, and enforcing mental representations of segregated social roles.⁹⁰ Stereotypical image outputs, in turn, contribute to confirming societal prejudices and feed into allocative harms, influencing the distribution of valuable social goods. The mitigation technique implemented by the developers of DALLE-2 seems to increase the diversity of population groups represented in image outputs. However, criticisms have been expressed towards the fact that diversity-related terms such as 'women' or 'black' were simply added to generic prompts to increase representativeness, thereby treating some of the symptoms of algorithmic bias without treating its root causes.⁹¹

88. AI modelling tool developed by Oxford academic incorporated into Amazon anti-bias software, Oxford Internet Institute, 21 April 2021, available at: <https://www.oii.ox.ac.uk/news/releases/ai-modelling-tool-developed-by-oxford-academics-incorporated-into-amazon-anti-bias-software-2/> (last accessed 29 September 2022)

89. See OpenAI, "Reducing Bias and Improving Safety in DALL-E 2" (18 July 2022) available at: <https://openai.com/blog/reducing-bias-and-improving-safety-in-dall-e-2/>.

90. Pamela Mishkin et al, "DALL-E 2 Preview - Risks and Limitations" (2022) available at: <https://github.com/openai/dalle-2-preview/blob/main/system-card.md#bias-and-representation>.

91. Matthew Sparkes, "AI art tool DALL-E 2 adds 'black' or 'female' to some image prompts", New Scientist (22 July 2022) available at: <https://www.newscientist.com/article/2329690-ai-art-tool-dall-e-2-adds-black-or-female-to-some-image-prompts/>.

While these are positive examples of existing governance efforts addressing algorithmic discrimination in the industry, it is important to highlight their limits.

The limits of technical solutions: debiasing and bias mitigation

First, **technical debiasing and bias mitigation solutions** cannot solve the problem of algorithmic discrimination in their own. As forcefully pointed out by Balayn and Gürses, “[d]ebiasing relies on conceptualisations of bias that do not capture the complexity of discrimination due to the limitations of the machine learning set-up.”⁹² Debiasing cannot redress algorithmic discrimination in a comprehensive or effective manner for two main reasons: On the one hand, these techniques focus exclusively on inputs and outputs of AI systems **without considering the context in which they are put to use**.⁹³ Debiasing techniques are algorithm-centric and **fail to consider the machine-human interaction** points that are also a source of bias.⁹⁴

On the other hand, **debiasing techniques themselves have not yet reached a development stage that allows for deployment across the board**: “[the] use cases are limited, the proposed conceptualisations of bias can oversimplify matters of discrimination, and the effectiveness and usability of debiasing methods and auditing tools are yet to be established”⁹⁵ The practical application of debiasing techniques is also a challenge because of difficulties surrounding the access to sensitive data as well as contextual variations across use cases.⁹⁶ For instance, anti-discrimination law might require different conceptions of fairness to intervene across different use cases or at different stages of the same use case, which are difficult to translate into technical metrics as well as difficult to reconcile with each other.

This leads to the **question of what it means for an algorithm to be ‘fair’?** A vast amount of research in computer science is dedicated to algorithmic ‘fairness’. Fairness approaches are sometimes presented as being able to ensure the ethical and legal compliance of algorithmic systems. Yet, **‘bias’ and ‘fairness’ are technical notions that do not neatly overlap with their ethical and legal counterparts**. In discrimination law, in particular, the prohibition

92. Balayn A and Gürses S, Beyond Debiasing: Regulating AI and its inequalities (European Digital Rights 2021), 51 available at: https://edri.org/wp-content/uploads/2021/09/EDRI_Beyond-Debiasing-Report_Online.pdf.

93. See *ibid*, 12, 64.

94. See *ibid*, 50.

95. *Ibid*, 12, 50.

96. See *ibid*.

on bias will be limited to those targeting or otherwise negatively impacting protected groups. Removing such biases at one point of the AI lifecycle might yield fairness from a technical perspective, nevertheless that might not adequately satisfy existing legal obligations pertaining to equality throughout the AI lifecycle.

In addition, computer scientists have developed a **wide range of definitions of fairness**, some of which are contradictory. Hence, **depending on the definition, an algorithm might be technically fair without necessarily complying with anti-discrimination law.**⁹⁷ From a mathematical standpoint, there are several ways to achieve a fair outcome, and they all relate to different perceptions and interpretations of fairness itself. For example, conceptualisations of fairness range from giving everyone the “same opportunity” while ignoring their wildly different starting points, to recognising the differences between people and giving some individuals a temporary advantage to counterbalance a disadvantage.⁹⁸ It could be argued for example, that treating a minority applicant the “same” when it comes to the provision of a loan may be fair. However, if due to historic and entrenched racism, that minority group has a higher risk of losing a job and thus being unable to repay the loan through no fault of their own, the application of fairness as simply the equalisation of outputs may lead to further entrenchment of inequality as those applicants may see their credit ratings further reduced.

Definitions of ‘fairness as accuracy’ and debiasing techniques aiming to acquire *more* data and building *more* accurate algorithmic systems also present important limits. While so-called “**accuracy-affecting injustices**” stemming from issues pertaining to data representativeness, data collection and data processing practices can be resolved via changes to data policies aiming to increase accuracy in algorithmic decision-making,⁹⁹ biases resulting from past injustices require different types of solutions. So-called “**nonaccuracy-affecting injustices**” give rise to data biases that cannot be addressed

97. See the discussion around differing ways of measuring bias and diverge definitions of fairness in the example of the COMPAS recidivism risk prediction system: Angwin, Julia, et al. “Machine bias.” *Ethics of Data and Analytics*. Auerbach Publications, 2016. 254-264 and Jeff Larson, Surya Mattu, Lauren Kirchner and Julia Angwin, “How We Analyzed the COMPAS Recidivism Algorithm” (2016) ProPublica available at: <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>.

98. With a view to substantive and transformative equality, so-called temporary special measures or positive action provide special support or a provisional advantage to a disadvantaged group so as to transform an unequal status quo in the long-term. See the discussion in section 3 of this study.

99. Hellman, Deborah. “Big Data and Compounding Injustice.” *Journal of Moral Philosophy*, forthcoming, Virginia Public Law and Legal Theory Research Paper 2021-27 (2021).

via improvements in data collection practices.¹⁰⁰ They reflect facts that are accurate but problematic because they result from historical discrimination and exclusion. **Only policies targeting the root causes and effects of such inequality can redress this type of bias.** For example, if an HR service wanted to automatise recruitment by predicting which candidates would be top performers, integrating more data about past recruitments will not address the causes of gender bias, which lie in gender segregation on the labour market, glass ceiling issues, the gender pay gap, gender stereotypes, etc.

Because of these limitations, solutionist narratives of debiasing should be debunked. **If at all, debiasing can only be one element of a broader anti-discrimination strategy** in relation to algorithmic systems. **Such a strategy should centre on human rights and socio-legal intervention as well as taking into account the whole deployment cycle of algorithmic decision-making systems** ranging from the formulation of the problem to address, to the context of implementation of the system, its actual performance and its practical impact. In addition, as pointed out by Balayn and Gürses, AI service providers should not enjoy wide discretion in choosing the strategies to prevent the discriminatory impact of their systems.¹⁰¹ Rather, **democratic control and regulatory safeguards should establish a framework around accepted fairness and anti-discrimination approaches, taking full account of technical limitations and of the need to address the root causes of algorithmic discrimination.** The participation of end-users directly affected by these systems, and in particular minority groups, should also be ensured. As highlighted in our recommendations below, this should also apply to standard-setting activities.

The limits of bias audits: access to data and diverging standards

Second, **auditing biases** has been presented as another potential solution to address algorithmic discrimination. Yet, problems arise in relation to access to data and diverging standards.

Auditing is defined as “a range of approaches to review algorithmic processing systems” which “can take different forms, from checking governance documentation, to testing an algorithm’s outputs, to inspecting its inner

100. Ibid.

101. See *ibid.*, 11.

workings”.¹⁰² It has been suggested that auditing could be used as a preventive safeguard against the release of discriminatory algorithmic systems on the market.¹⁰³ However, the **lack of access to equality data**, **GDPR-related uncertainties** on the permitted processing of sensitive categories of data and **the lack of uniformly accepted standards** makes auditing algorithms for discrimination challenging.

On the one hand, legal scholars are **uncertain about whether the GDPR allows processing of sensitive categories of personal data for debiasing** or more broadly for **anti-discrimination purposes**.¹⁰⁴ On the other hand, the **lack of equality data**, stemming from often restrictive equality data collection practices in Europe, raises issues when it comes to identifying inequality in specific domains such as access to housing, education, healthcare, work, etc. for various protected groups of population.¹⁰⁵ It limits access to accurate information about ground truth and the extent of structural inequality in society.

This problem of accessing sensitive data should also be considered in the broader context of **data extraction and exploitation** by big tech firms. **Access to such data for discrimination auditing and anti-discrimination purposes in general should therefore be entrusted to other entities, possibly including equality bodies, labour inspectorates, CSOs with a legitimate interest** in the sense of Art. 11 and 12 and Art. 13 and 14 of the EU equality directives 2000/43/EC and 2000/78/EC, etc. The development of **more systematic, ethical and regulated equality data collection** throughout Europe, which

102. Digital Regulation Co-operation Forum, “Auditing algorithms: the existing landscape, role of regulators and future outlook” (2022) available at: <https://www.gov.uk/government/publications/findings-from-the-drcf-algorithmic-processing-workstream-spring-2022/auditing-algorithms-the-existing-landscape-role-of-regulators-and-future-outlook>.

103. See Kim PT, ‘Auditing Algorithms for Discrimination’ (2017) 166 University of Pennsylvania Law Review Online 189.

104. Van Bekkum, Marvin and Zuiderveen Borgesius, Frederik, Using Sensitive Data to Prevent Discrimination by AI: Does the GDPR Need a New Exception? (2022) available at: <http://dx.doi.org/10.2139/ssrn.4104823> (last accessed 28 July 2022).

105. See European Commission, Analysis and comparative review of equality data collection practices in the European Union : legal framework and practice in the EU Member States (Publications Office, 2017) available at: <https://data.europa.eu/doi/10.2838/6934> (last accessed 28 July 2022); Lilla Farkas, Analysis and comparative review of equality data collection practices in the European Union : data collection in the field of ethnicity (Publications Office, 2020) available at: <https://data.europa.eu/doi/10.2838/447194> (last accessed 28 July 2022); Ringelheim, Julie, “Processing Data on Racial or Ethnic Origin for Antidiscrimination Policies: How to Reconcile the Promotion of Equality with the Right to Privacy?” (2007) NYU School of Law Jean Monnet Working Paper No. 08/06, available at: <http://dx.doi.org/10.2139/ssrn.983685> (last accessed 28 July 2022).

numerous actors in the field, including equality bodies, have long been advocating for, would also be progress for algorithmic auditing purposes. Inspiration could come from the UK, where, as part of the “Data: a new direction strategy” policy,¹⁰⁶ the Government has announced that a new condition will be introduced under the Data Protection Act (DPA) 2018 to allow for the processing of special category data for the monitoring and mitigation of algorithmic bias.

Furthermore, there are neither legal obligations nor uniform standards for algorithmic auditing yet. Various methodologies have been proposed.¹⁰⁷ Some of the toolkits developed by researchers have been adopted by major companies, for instance the ‘Aequitas’ instrument developed at the Oxford Internet Institute and adopted by Amazon.¹⁰⁸ Nonetheless, developing uniform regulatory standards for algorithmic auditing in the field of non-discrimination would substantially increase legal certainty for providers. This would also foster public trust in algorithmic systems. Finally, uniform regulatory standards for algorithmic auditing would enhance companies’ take up of discrimination audits, which would in turn provide useful information for potential victims to assess the opportunity of taking (legal) action and comprehensible evidentiary material to judges.

6) Representation and participation issues: The lack of diversity and inclusion in the AI industry

The under-representation of disadvantaged groups in professional communities involved with the development of AI is an important dimension of the problem of algorithmic discrimination. The lack of diversity and inclusion in these communities means that women and under-represented groups do not (sufficiently) participate in the crafting of algorithmic technologies, with the consequence that **they cater suboptimally to the needs of these groups, disadvantages them or even erases them entirely.** A

106. Data: a new direction - government response to consultation, 22 June 2022, available at: <https://www.gov.uk/government/consultations/data-a-new-direction/outcome/data-a-new-direction-government-response-to-consultation> (last accessed: 28 July 2022)

107. For a review, see e.g., Jack Bandy, (2021) ‘Problematic Machine Behaviour: A Systematic Literature Review of Algorithm Audits.’ Forthcoming, Proceedings of the ACM (PACM) Human-Computer Interaction, CSCW ’21.

108. See Saleiro, P, Kuester, B, Hinkson, L, London, J, Stevens, A, Anisfield, A, Rodolfa, KT, Ghani, R (2018) ‘Aequitas: A Bias and Fairness Audit Toolkit.’ Arxiv and Oxford Internet Institute (2021) ‘AI modelling tool developed by Oxford Academics incorporated into Amazon anti-bias software’ available at: <https://www.oii.ox.ac.uk/news-events/news/ai-modelling-tool-developed-by-oxford-academics-incorporated-into-amazon-anti-bias-software-2/>.

survey issued by the Council of Europe for the purpose of the present Study shows that **most responding State Parties to the ECHR are aware of the diversity issue** in the AI industry. **State Parties highlight the need to steer more women and minority groups towards STEM (science, technology, engineering, mathematics) disciplines** as this is perceived as a major factor contributing to discriminatory AI.

Some notable examples of AI bias due to lack of diversity have been exposed in a report by the AI Now Institute, founded by ex-Google executive Meredith Whittaker and principal researcher at Microsoft Research Kate Crawford.¹⁰⁹ These include image recognition services which classified black people as gorillas and Amazon technology failing to recognise users with darker skin colours. The thesis of the report (reflecting a widely held view in the broader academic, policy and AI community) is that examples such as these occur due to “blind spots” because developers design and test models based on their own standpoint. The lack of a diverse workforce leads to a limited perspective and can result in bias that may be difficult to detect and correct before it leads to discrimination.

In addition to the widespread problem of **implicit bias**, a homogenous group is likely to have a **truncated outlook influenced by similar identities and experiences**. As an example, the Google AI Experiments programme developed a game called “Quick, Draw!” In the game, people were asked to draw pictures of everyday things like shoes to train a model.¹¹⁰ All five of the game’s developers at Google were men. They and early users of the game drew men’s sneakers to represent a shoe. This resulted in a game which did not know that high heels were also shoes. This was not an intentional error; it was simply shaped by the perspective of the dominant representative group designing algorithms in the technology industry. **As such, any algorithm built by a majority group is at risk of failing to embed perspectives of marginalised minority groups, resulting in algorithms that only work for the majority.**

Diversity matters as it provides holistic approaches in making AI technologies more responsible. It helps address challenges faster and clearer as local knowledge and front-line experience will be embedded in the core of every decision-making or working process. Getting the right mix of minds in the

109. Sarah Myers West, Meredith Whittaker and Kate Crawford, *Discriminating Systems: Gender, Race, and Power in AI*, AI Now Institute NYU, April 2019, available at: <https://ainowinstitute.org/discriminatingystems.pdf> (last accessed: 27 July 2022).

110. Josh Lovejoy, *Fair Is Not the Default – Why building inclusive tech takes more than good intentions*, 15 February 2018, <https://design.google/library/fair-not-default/> (last accessed: 28 July 2022).

room is essential to gain the necessary insight to address bias and gain competitive advantage. **Diversity should therefore be viewed as being “mission critical” when it comes to innovation.** This should translate in more diverse recruitment policies in educational and professional communities involved with the development and use of AI systems. As argued in section 3, legal obligations revolving around the notion of positive action could play a major role in this regard. In addition, diversity policies in educational and professional recruitment should be complemented by adequate training.

The *AI Now* (New York University) report¹¹¹ identified a “diversity crisis” in the AI sector, especially in the global technology industry, which is overwhelmingly white and male, and asserts that this has contributed to algorithmic gender and racial biases. A 2020 World Economic Forum report¹¹² painted a similarly grim picture: despite talk of greater inclusion, women’s representation in tech-related jobs has declined by 32% since 1990. According to a study launched by the EU Commission in 2016, “only 24 out of every 1000 female graduates had an ICT related subject in her portfolio”. When it comes to employment, only 6 of those girls and women finally found a job in the digital sector.¹¹³

A Canadian start-up found that women make only 12% of leading machine learning researchers.¹¹⁴ Another report¹¹⁵ by New York University – *Discriminating Systems – Gender, Race, and Power in AI* asserts discrimination in AI systems was associated with the lack of diversity in the teams that work these technologies. **Whether the focus is on mitigation of bias in input processes, or fairness in outcomes, diversity and inclusion is one of the most powerful tool companies have at their disposal.** The blind spots created by the

111. Kari Paul, ‘Disastrous’ lack of diversity in AI industry perpetuates bias, study finds, *The Guardian*, 17 April 2019, available at: <https://www.theguardian.com/technology/2019/apr/16/artificial-intelligence-lack-diversity-new-york-university-study> (last accessed: 27 July 2022).

112. Ronit Avi and Rana El Kaliouby, Here’s why AI needs a more diverse workforce, *World Economic Forum*, 21 September 2020 <https://www.weforum.org/agenda/2020/09/ai-needs-diverse-workforce/> (last accessed: 27 July 2022).

113. Women in AI: Promoting inclusive participation across society, Aimee Van WYNSBERGH, European AI Alliance, available at: <https://futurium.ec.europa.eu/en/european-ai-alliance/blog/women-ai-promoting-inclusive-participation-across-society?language=hu> (last accessed: 31 August 2022).

114. Archie de Berker, Women in Machine Learning: Negar Rostamzadeh, 20 February 2018, available at: <https://medium.com/element-ai-research-lab/women-in-machine-learning-negar-rostamzadeh-dbb58dc75e81> (last accessed: 31 August 2022).

115. Sarah Myers West, Meredith Whittaker and Kate Crawford, *Discriminating Systems: Gender, Race, and Power in AI*, AI Now Institute NYU, April 2019, available at: <https://ainowinstitute.org/discriminatingystems.pdf> (last accessed: 27 July 2022).

lack of diversity – diversity of education, perspectives, life experiences and backgrounds – make it more challenging to anticipate biases in algorithmic systems and their potential impact on different individuals and groups.

Already marginalised groups are systematically and disproportionately put at more risk of being harmed by algorithmic decision-making tools that do not represent their perspectives and interests. Beyond the moral imperative of preventing systemic racial and gender discrimination in designing new AI tools, there is also an economic one. Research has demonstrated that “companies in the top quartile for gender diversity have been 21% more likely to experience above-average profitability, while ethnic and cultural diversity correlates with a 33% increase in performance.”¹¹⁶

116. The five business benefits of a diverse team, CMI, 3 July 2019, available at: <https://www.managers.org.uk/knowledge-and-insights/listicle/the-five-business-benefits-of-a-diverse-team/> (last accessed: 31 August 2022).

Section 2

The legal and policy landscape in Europe: strengths and shortcomings

There is **general awareness among policy makers** that, alongside opportunities, AI brings the risks of solidifying and perpetuating existing inequalities. In a survey issued by the Council of Europe to gauge the views of the members and observers of the GEC and the CDADI, more than **80% of respondents viewed AI as posing risks to human rights. 40% of respondents identified a direct risk of gender discrimination.**

Several initiatives are taking place across governments, and they encompass several issues, from female participation in STEM fields, to deepfakes and cyberbullying and algorithmic discrimination. For example, some countries, like **Finland**, have addressed the issue of the lack of transparency in algorithmic systems leading to discrimination head on, issuing recommendations and guidance to raise awareness of the problem.¹¹⁷ The **Netherlands** has adopted a 'Fundamental rights and algorithms Impact Assessment' that includes a 'Non-discrimination by design guideline'.¹¹⁸ The Dutch Parliament has recently adopted a motion rendering human rights impact assessments compulsory for public institutions using algorithms.¹¹⁹

117. Automaattisessa päätöksenteossa on turvattava virkavastuu ja hyvän hallinnon toteutuminen, available at: <https://valtioneuvosto.fi/-/10623/automaattisessa-paatoksenteossa-on-turvattava-virkavastuu-ja-hyvan-hallinnon-toteutuminen> (last accessed: 28 July 2022).

118. Ministry of the Interior and Kingdom Relations, 'Fundamental rights and algorithms Impact Assessment' (March 2022) available at: <https://www.government.nl/binaries/government/documenten/reports/2022/03/31/impact-assessment-fundamental-rights-and-algorithms/Fundamental+Rights+and+Algorithms+Impact+Assessment.pdf>.

119. See European Center for Not-for-Profit Law, "Netherlands sets precedent for human rights safeguards in use of AI" (2022) available at: <https://ecnl.org/news/netherlands-sets-precedent-human-rights-safeguards-use-ai>.

The **Austrian** government has published an action plan on deepfakes, including diverse measures to tackle the problem. In **Finland**, Aurora AI aims to guide citizens, especially young people, to the services they need by means of artificial intelligence. If, as a result, young people find the services they need better, this is likely to promote equality, for example in access to services or in the provision of assistance and support. The **Portuguese** Agency for Administrative Modernisation (AMA) has developed - with the help of the Commission for Citizenship and Gender Equality and other relevant stakeholders - the “Guide for the use of Artificial Intelligence in Public Administration”. The guide is designed to address the concerns of non-discrimination in general and the protection of individual and collective rights in the development of algorithmic systems. It draws attention to the reliability and representativeness of the data to be collected and processed, and emphasises the issues associated with ethics, justice, transparency, accountability and understanding of the systems.

Yet, national responses are largely **uncoordinated**. While legislators such as the European Union are in the process of adopting a uniform regulatory framework on AI, **the Council of Europe could exert wide-ranging regulatory influence in the field of human rights**. Where the EU is advocating for a ‘human-centric AI’, regulatory action by the Council of Europe could foster a distinct **human-rights-based approach to AI**.

This section of the Study highlights **which existing legal instruments at Council of Europe level can be used to address various dimensions of the problem of algorithmic discrimination**, ranging from **non-discrimination to data protection and privacy law to sectoral regulations**. It also briefly maps existing and forthcoming EU legal instruments and shows that both frameworks present **shortcomings and uncertainties when it comes to addressing algorithmic discrimination**. These gaps call for regulatory action at Council of Europe level, some possible contours of which are highlighted in Section 3.

I. Discrimination and equality: legal and policy instruments and their limits

This section highlights the existing legal instruments that provide a legal basis for combatting algorithmic discrimination and related forms of algorithmic violence.

1) Binding legal instruments of the Council of Europe

The European Convention on Human Rights

Article 14 ECHR and Art. 1 of Protocol No. 12 lay out a prohibition on discrimination that provides a **legal basis for banning algorithmic discrimination**.

Article 14 ECHR prohibits discrimination based on an open-ended list of protected characteristics:

“The enjoyment of the rights and freedoms set forth in [the] Convention shall be secured without discrimination on any ground such as sex, race, colour, language, religion, political or other opinion, national or social origin, association with a national minority, property, birth or other status.”

Article 14 ECHR is not a self-standing provision, meaning that it can only be invoked in association with a claim that another substantive right protected by the ECHR has been violated.

Entered into force in 2005, **Protocol No. 12 to the Convention** has so far been ratified by 20 out of 46 state parties to the ECHR. Article 1 lays out a free-standing general prohibition of discrimination:

“1. The enjoyment of any right set forth by law shall be secured without discrimination on any ground such as sex, race, colour, language, religion, political or other opinion, national or social origin, association with a national minority, property, birth or other status.

2. No one shall be discriminated against by any public authority on any ground such as those mentioned in paragraph 1.”

The Istanbul Convention

The Council of Europe Convention on preventing and combating violence against women and domestic violence (the Istanbul Convention) provides a **legal basis for prohibiting digital violence against women, including algorithmic stereotyping, and online violence such as cyber-harassment, bullying and online sexist hate speech**.

The Istanbul Convention was adopted in 2011, entered into force in 2014 and has been ratified by 37 state parties. It recognises gender-based violence (GBV) against women as a form of discrimination. Its provisions focus on prevention, protection, prosecution and the development of integrated policies in relation to combating violence against women and domestic violence.

The Istanbul Convention provides a legal basis for **addressing the continuum of online and offline violence against women**.¹²⁰ It requires that stalking, sexual harassment, and psychological violence, including when committed via information and communication technology (ICT), are sanctioned. Moreover, it gives a clear mandate to the public authorities of state parties to address the societal roots and the online embodiments of gender-based violence. Particularly relevant to issues of online gender-based violence is also **Art. 17** on “participation of the private sector and the media” which states that:

“Parties shall encourage the private sector, the information and communication technology sector and the media, with due respect for freedom of expression and their independence, to participate in the elaboration and implementation of policies and to set guidelines and self-regulatory standards to prevent violence against women and to enhance respect for their dignity”.

The Framework Convention for the Protection of National Minorities

The Framework Convention for the Protection of National Minorities provides a **legal basis for combatting algorithmic discrimination on grounds of national minority status as well as online violence such as hate speech**. Entered into force in 1998, the Convention counts 39 state parties. In its **Art. 4**, the Convention states that:

“1. The Parties undertake to guarantee to persons belonging to national minorities the right of equality before the law and of equal protection of the law. In this respect, any discrimination based on belonging to a national minority shall be prohibited.

2. The Parties undertake to adopt, where necessary, adequate measures in order to promote, in all areas of economic, social, political and cultural life, full and effective equality between persons belonging to a national minority and those belonging to the majority. In this respect, they shall take due account of the specific conditions of the persons belonging to national minorities.”

Art 6(2) lays out that *“The Parties undertake to take appropriate measures to protect persons who may be subject to threats or acts of discrimination, hostility or violence as a result of their ethnic, cultural, linguistic or religious identity.”*

120. See the General Recommendation No. 1 on the Digital Dimension of Violence against Women adopted by the Group of Experts on Action against Violence against Women and Domestic Violence (GREVIO), available at: <https://www.coe.int/en/web/istanbul-convention/general-recommendation>.

Art 9 relating to freedom of expression, which states that *“The Parties shall ensure, within the framework of their legal systems, that persons belonging to a national minority are not discriminated against in their access to the media”*, could become particularly relevant for issues of discrimination on social media platforms, cyberharassment and hate speech.

The European Charter for Regional or Minority Languages

Entered into force in 1998, the Charter has been ratified by 25 countries so far. **Art 7(2)** of the Charter lays out that *“The Parties undertake to eliminate, if they have not yet done so, any unjustified distinction, exclusion, restriction or preference relating to the use of a regional or minority language and intended to discourage or endanger the maintenance or development of it”*. Again, in principle **this provision extends to the algorithmic and online realms, where it can be relied on to address digital discrimination in its many forms**. In addition, the Committee of Experts of the European Charter for Regional or Minority Languages recently published a statement highlighting how ‘AI applications may facilitate the daily use of regional or minority languages and support authorities in promoting them in accordance with the Charter’ and ‘encourag[ing] states to promote the inclusion of regional or minority languages into research and study on AI’.¹²¹

The European Social Charter

To date 43 members of the Council of Europe have ratified either the European Social Charter (ETS No. 35), adopted in 1961, or the Revised European Social Charter (ETS No.163), adopted in 1996. In the revised Charter, **Art E** on ‘**Non-discrimination**’ provides that **‘the enjoyment of the rights set forth in this Charter shall be secured without discrimination on any ground such as race, colour, sex, language, religion, political or other opinion, national extraction or social origin, health, association with a national minority, birth or other status’**. In addition, Art. 20 guarantees ‘the right to equal opportunities and equal treatment in matters of employment and occupation without discrimination on the grounds of sex’. Both provisions are relevant to discrimination induced by algorithmic systems, especially as the European Social Charter and its revised version focus on fundamental social rights that relate to employment and working conditions, housing, education, health, medical assistance and social protection, i.e. areas which have been deeply impacted by new forms of algorithmic management.

121. Statement of the Committee of Experts of the European Charter for Regional or Minority Languages on the promotion of regional or minority languages through artificial intelligence (2022) available at: <https://rm.coe.int/declaration-ai-en/1680a657ff>.

2) Relevant policy instruments of the Council of Europe

A number of **non-binding standards and policy instruments** complement the binding legal provisions and are **relevant when it comes to addressing the discriminatory effects of AI and algorithmic decision-making**.

In March 2019, the “**Recommendation on Preventing and Combating Sexism**” drafted by the Gender Equality Commission was adopted by the Council of Ministers.¹²² It recognises that “[t]he internet has provided a new dimension for the expression and transmission of sexism, especially of sexist hate speech, to a large audience, even though the roots of sexism do not lie in technology but in persistent gender inequalities”.¹²³ It enjoins member states to “integrate a gender equality perspective in all policies, programmes and research in relation to artificial intelligence to avoid the potential risks of technology perpetuating sexism and gender stereotypes”.¹²⁴ The recommendation also foresees a positive role for AI as it requires State Parties to “examine how artificial intelligence could help to close gender gaps and eliminate sexism”.¹²⁵ It lists key aspects such as women’s and girls’ participation in IT education and industries, the mainstreaming of gender equality in the design of data-driven instruments, awareness-raising as regards gender bias in big data, transparency and accountability. In turn, the recent recommendation “**On combating hate speech**” co-drafted by the Steering Committee on Anti-Discrimination, Diversity and Inclusion (CDADI) and the Steering Committee on Media and Information Society (CDMSI) indicates that “internet intermediaries should identify expressions of hate speech that are disseminated through their systems and act upon them in the framework of their corporate responsibility”.¹²⁶

The **Council of Europe Gender Equality Strategy 2018-2023** also recognises that “sexism and discrimination against women includ[e] **sexist hate speech online**” as well as online gender-based violence.¹²⁷ In addition, in 2021, the Council of Europe Group of Experts on Action against Violence against Women and Domestic Violence (GREVIO), which monitors the implementation of the

122. Council of Europe, “Preventing and combating sexism”, Recommendation CM/Rec(2019)1 adopted by the Committee of Ministers of the Council of Europe (27 March 2019), available at <https://rm.coe.int/prems-055519-gbr-2573-cmrec-2019-1-web-a5/168093e08c>.

123. Ibid.

124. Recommendation II.B.7, *ibid.*, p. 19.

125. Ibid.

126. Council of Europe, Recommendation CM/Rec(2022)16[1] of the Committee of Ministers to member States on combating hate speech (20 May 2022), [30].

127. Council of Europe Gender Equality Strategy 2018-2023 adopted by the Committee of Ministers (March 2018), p. 10, 16, 18, available at <https://rm.coe.int/prems-093618-gbr-gender-equality-strategy-2023-web-a5/16808b47e1>.

Istanbul Convention, adopted its General Recommendation No.1 on the **digital dimension of violence against women**, which highlights legal issues around online sexual harassment, stalking and the digital dimension of psychological violence.¹²⁸

In May 2022, the Committee of Ministers adopted a new **Recommendation on combating hate speech** jointly drafted by the Steering Committees on Anti-Discrimination, Diversity and Inclusion (CDADI) and on Media and Information Society (CDMSI),¹²⁹ It recognises the existence of a “**power asymmetry between some digital platforms and their users**” and makes recommendations for tackling **online hate speech in relation to policies pertaining to content moderation, micro-targeting and online advertising, content amplification, recommender systems and underlying data collection strategies**.

In May 2022, the Committee of Ministers adopted a “**Recommendation on protecting the rights of migrant, refugee and asylum seeking women and girls**” which demands that **human rights impact assessments are conducted before AI and automated decision making systems are introduced in the field of migration and that the design, development and application of such systems are non-discriminatory**.¹³⁰ It also calls for involving refugee, asylum-seeking and migrant women and representative CSOs “in discussions on the development and deployment of new technologies affecting them”.

Other instruments such as the **Guidelines of the Committee of Ministers of the Council of Europe on upholding equality and protecting against discrimination and hate during the Covid-19 pandemic and similar crises in the future** mention the need to ensure that “digital tools for dealing with the crisis and the resulting risks” “are not discriminatory against persons belonging to vulnerable groups or otherwise violate their rights”¹³¹

128. Group of Experts on Action against Violence against Women and Domestic Violence, General Recommendation No. 1 on the digital dimension of violence against women (20 October 2021) available at: <https://rm.coe.int/grevio-rec-no-on-digital-violence-against-women/1680a49147> (last accessed 22 July 2022).

129. Council of Europe, Recommendation CM/Rec(2022)16 on combating hate speech, available at https://search.coe.int/cm/Pages/result_details.aspx?ObjectId=0900001680a67955#_ftn1.

130. Council of Europe, Recommendation CM/Rec(2022)17 of the Committee of Ministers to member States on protecting the rights of migrant, refugee and asylum-seeking women and girls, [22]-[25] available at https://search.coe.int/cm/Pages/result_details.aspx?ObjectId=0900001680a69407.

131. Steering Committee on Anti-Discrimination, Diversity And Inclusion (CDADI), Guidelines of the Committee of Ministers of the Council of Europe on upholding equality and protecting against discrimination and hate during the Covid-19 pandemic and similar crises in the future (2020), [27] available at: <https://rm.coe.int/prems-066521-gbr-2530-cdadi-guidelines-web-a5-corrige/1680a3d50c>.

Together, these recommendations address a number of issues contributing to algorithmic discrimination as pointed out earlier in this Study: the lack of diversity, equal representation and equal participation in educational and professional fields related to the AI industry, the lack of binding obligation to mainstream equality-related concerns in the development of algorithmic systems and the lack of clearly defined accountability mechanisms.

Furthermore, the Council of Europe has adopted specific policy instruments relating to human rights in the digital space. In 2020, the Council of Ministers adopted **Recommendation CM/Rec(2020)1 on the human rights impacts of algorithmic systems**. It pays special attention to discrimination by requiring, for example, that ‘private sector actors that design, develop or implement algorithmic systems [...] follow a standard framework for human rights due diligence to avoid fostering or entrenching discrimination throughout all life-cycles of their systems’ and that ‘[t]hey seek to ensure that the design, development and ongoing deployment of their algorithmic systems do not have direct or indirect discriminatory effects on individuals or groups that are affected by these systems, including on those who have special needs or disabilities or who may face structural inequalities in their access to human rights.’¹³² Other relevant policy instruments include **Recommendation CM/Rec(2022)13 of the Committee of Ministers to member States on the impacts of digital technologies on freedom of expression** and various sets of guidelines on facial recognition,¹³³ content moderation¹³⁴ and artificial intelligence and data protection.¹³⁵

132. Committee of Ministers, Recommendation CM/Rec(2020)1 of the Committee of Ministers to member States on the human rights impacts of algorithmic systems (Adopted by the Committee of Ministers on 8 April 2020 at the 1373rd meeting of the Ministers’ Deputies), available at: https://search.coe.int/cm/pages/result_details.aspx?objectid=09000016809e1154.

133. Consultative Committee of the Convention for the protection of individuals with regard to automatic processing of personal data (Convention 108), Guidelines on facial recognition (2021) available at: <https://edoc.coe.int/en/artificial-intelligence/9753-guidelines-on-facial-recognition.html>.

134. Council of Europe Guidance Note: content moderation. Best practices towards effective legal and procedural frameworks for self-regulatory and co-regulatory mechanisms of content moderation (adopted by the Steering Committee for Media and Information Society (CDMSI)) (2021) available at: <https://rm.coe.int/content-moderation-en/1680a2cc18>.

135. Consultative Committee of the Convention for the protection of individuals with regard to automatic processing of personal data (Convention 108), Guidelines on artificial intelligence and data protection (2019) available at: <https://rm.coe.int/guidelines-on-artificial-intelligence-and-data-protection/168091f9d8>.

3) Comparative insights: other relevant European and international provisions

The European Union also has a very developed legal framework on discrimination and equality. **Art 21(1)** of the **EU Charter of Fundamental Rights** prohibits discrimination “*on any ground such as sex, race, colour, ethnic or social origin, genetic features, language, religion or belief, political or any other opinion, membership of a national minority, property, birth, disability, age or sexual orientation*” and **Art. 21(2)** states that “[w]ithin the scope of application of the Treaties and without prejudice to any of their specific provisions, any discrimination on grounds of nationality shall be prohibited”. **Art. 23** indicates that “[e]quality between women and men must be ensured in all areas, including employment, work and pay” and allows positive action. In secondary law, **Directive 2000/43/EC** guarantees equality on grounds of race or ethnic origin at work, in the access to goods and services and in education. **Directive 2000/78/EC** prohibits discrimination on grounds of disability, sexual orientation, religion or belief and age in the workplace and vocational training. **Directive 2004/113/EC** guarantees gender equality in the access to goods and services and so does **Directive 2006/54/EC** in relation to work.

In 2022, the European Commission published a “**European Declaration on Digital Rights and Principles for the Digital Decade**” that reflects the Commission’s wish to develop a “**human-centred AI**” and exposes the EU’s approach to digital transformation. The Commission’s rationale is that digital rights should ensure that EU citizens have access to digital technologies and are protected from their harmful consequences. Chapter III of the Declaration includes a commitment to “ensuring that algorithmic systems are based on suitable datasets to avoid unlawful discrimination and enable human supervision of outcomes affecting people”.¹³⁶

At United Nations level, a number of instruments protect against discrimination beyond existing general human rights instruments: in particular the **International Convention on the Elimination of All Forms of Racial Discrimination (CERD)**, the **Convention on the Elimination of All Forms of Discrimination against Women (CEDAW)**, and the **Convention on the Rights of Persons with Disabilities (CRPD)**. More specifically, the CERD Committee issued a **General recommendation No. 36 on preventing and combating racial profiling by law enforcement officials** in 2020. This document

136. European Commission, “European Declaration on Digital Rights and Principles for the Digital Decade” COM(2022) 28 final (Brussels 2022).

recognises how the use of artificial intelligence leads to entrenching racial inequalities and makes recommendations to prevent and redress racial bias and discrimination.

Although these legal and policy instruments do not stop at the borders of the digital world, their applicability to the various forms of algorithmic discrimination suffers a number of shortcomings.

4) Limits and uncertainties: where does algorithmic discrimination fall into the cracks?

This legal and policy patchwork addresses some of the discriminatory risks of AI and automated decision-making. Yet, **many uncertainties remain concerning the extent to which existing legal provisions can be used to promote equality and counter discrimination arising from the use of these technologies.** Hence, the aim of this subsection is to explore existing gaps in the equality and non-discrimination framework described above when it comes to algorithmic discrimination. **Three main issues arise:** (1) **the lack of neat overlap between existing concepts of direct and indirect discrimination and forms of algorithmic discrimination;** (2) **procedural issues linked to evidence and responsibility;** (3) **challenges linked to the protection of specific characteristics by the law.** As explained in Section 3, addressing those gaps calls for enforcing existing positive obligations to promote equality and mainstreaming preventive approaches to algorithmic discrimination under Art. 14 ECHR.

Qualification issues: direct vs indirect algorithmic discrimination

Although Article 14 ECHR does not distinguish between direct and indirect discrimination, the European Court of Human Rights (the Court) carved out the distinction in its case law.¹³⁷ **Direct discrimination** arises from “**a difference in the treatment of persons in analogous, or relevantly similar, situations**” and where this difference is “**based on an identifiable characteristic**”

137. This has been done by reference to EU equality law and the case law of the European Court of Justice, see *D.H. and Others v. The Czech Republic* Application no. 57325/00 (European Court of Human Rights, Grand Chamber, 13 November 2007), [184].

or “status”.¹³⁸ For example, where two workers are similarly qualified for a promotion but one is preferred over the other “because of” their sex, this would give rise to direct sex discrimination.

At the beginning of the 2000s, the **Court recognised the existence of indirect discrimination** where states “fail to treat differently persons whose situations are significantly different”.¹³⁹ It ruled in *DH* that “a difference in treatment may take the form of disproportionately prejudicial effects of a general policy or measure which, though couched in neutral terms, discriminates against a group”.¹⁴⁰ For instance, a neutrally formulated policy that would make the recruitment of candidates conditional on a minimum height might have indirectly discriminatory effects on women, who are on average smaller than men.

Once a *prima facie* finding of direct or indirect discrimination has been established, an **open justification system** applies whereby **discrimination** can only be found where there is “no objective and reasonable justification”.¹⁴¹ In other terms, both direct and indirect discrimination can be justified if it pursues a **legitimate aim** and if there is a “relationship of **proportionality between the means employed and the aim sought** to be realized”.¹⁴² Because the same justification regime applies in principle under both frameworks, qualifying algorithmic discrimination as direct or indirect has less significant repercussions on available means of redress than under EU law, where this qualification conditions the applicability of a closed or an open

138. See e.g., *Kjeldsen, Busk Madsen and Pedersen v. Denmark* Application no. 5095/71, 5920/72, 5926/72 (European Court of Human Rights, 7 December 1976), [56]; *Burden v. the United Kingdom* Application 13378/05 (European Court of Human Rights, Grand Chamber, 29 April 2008), [60]; *Carson and Others v United Kingdom* Application no. 42184/05 (European Court of Human Rights, Grand Chamber, 16 March 2010), [61], and more recently *Biao v. Denmark* Application no. 38590/10 (European Court of Human Rights, Grand Chamber, 24 May 2016), [89]. See also European Union Agency for Fundamental Rights and Council of Europe, *Handbook on European non-discrimination law* (Publications Office of the European Union 2018), 43 and European Court of Human Rights, *Guide on Article 14 of the European Convention on Human Rights and on Article 1 of Protocol No. 12 to the Convention* (Council of Europe 2020), 11.

139. *Thlimmenos v. Greece* Application no. 34369/97 (European Court of Human Rights, 2 April 2000), [44].

140. *D.H. and Others v. The Czech Republic* Application no. 57325/00 (European Court of Human Rights, Grand Chamber, 13 November 2007), [184].

141. Case “relating to certain aspects of the laws on the use of languages in education in Belgium” v. Belgium Application no 1474/62; 1677/62; 1691/62; 1769/63; 1994/63; 2126/64 (European Court of Human Rights, 23 July 1968), [10] at 34.

142. *Ibid*, see also *Marckx v. Belgium* Application no. 6833/74 (European Court of Human Rights, 13 June 1979), [33].

regime of justifications¹⁴³ Nonetheless, it is important to understand how courts, including the Court, will qualify algorithmic discrimination.

So far, it has been argued that algorithmic discrimination mainly falls within the framework of indirect discrimination, in particular because developers are unlikely to input protected characteristics in the datasets used to train algorithmic decision-making (ADM) systems.¹⁴⁴ According to Hacker, for example, “in machine learning contexts, indirect discrimination is the most relevant type of discrimination” while “[d]irect discrimination will be rare in algorithmic decision making, and largely limited to cases of implicit bias in labelling”.¹⁴⁵ Borgesius and Kelly-Lyth also respectively argue that “non-discrimination law prohibits many discriminatory effects of algorithmic decision-making, in particular through the concept of indirect discrimination”¹⁴⁶ and that “most biased algorithms will fall under the indirect discrimination framework”.¹⁴⁷

At least three arguments support this view: (1) Indirect discrimination captures situations where formally neutral measures produce disadvantage because they intervene in, and embed, an unequal social context.¹⁴⁸ This resonates with the ways in which data-driven technologies incorporate and perpetuate society’s unequal *status quo*.¹⁴⁹ (2) Indirect discrimination focuses on the structural dimension of discrimination.¹⁵⁰ This focus resonates with

143. Under EU law, direct discrimination cannot, in principle, be justified (safe closed exceptions), while indirect discrimination gives rise to a proportionality test with an open-ended regime of justifications.

144. This argument builds on an analogy with the US anti-discrimination framework, see eg Solon Barocas and Andrew D. Selbst, ‘Big Data’s Disparate Impact’ (2016) 104 California law review 671. Yet, the distinction between direct and indirect discrimination in ECHR law differs from the US distinction between notions of “disparate treatment” and “disparate impact”.

145. Hacker, ‘Teaching fairness to artificial intelligence: Existing and novel strategies against algorithmic discrimination under EU law’, 1152-1153.

146. Zuiderveen Borgesius, ‘Strengthening legal protection against discrimination by algorithms and artificial intelligence’, 1578. He nevertheless acknowledges a range of enforcement issues.

147. Aislinn Kelly-Lyth, ‘Challenging Biased Hiring Algorithms’ (2021) 41 Oxford Journal of Legal Studies 899, 906.

148. See Tobler, *Limits and potential of the concept of indirect discrimination*, 85. On the perpetrator’s vs. the victim’s perspective, see Alan David Freeman, ‘Legitimizing Racial Discrimination Through Antidiscrimination Law: A Critical Review of Supreme Court Doctrine’ (1978) 62 Minnesota Law Review 1049.

149. See Anna Lauren Hoffmann, ‘Where fairness fails: data, algorithms, and the limits of anti-discrimination discourse’ (2019) 22 Information, Communication & Society 900.

150. See Hugh Collins and Tarunabh Khaitan, ‘Indirect Discrimination Law: Controversies and Critical Questions’ in Hugh Collins and Tarunabh Khaitan (eds), *Foundations of Indirect Discrimination Law* (1 edn, Hart Publishing 2018), 19.

the fact that machine learning (ML) algorithms derive rules from group patterns. (3) The concept of indirect discrimination allows addressing distinctions not based on legally protected grounds that in effect impact protected groups.¹⁵¹ Since such proxy discrimination is one of the prevailing forms of algorithmic discrimination, as will be explained below, the framework of indirect discrimination presents a further advantage.

Despite the consensus on classifying algorithmic discrimination as indirect, **such a qualification “by default” raises a number of doctrinal and procedural issues.**¹⁵² As recent research shows, the notion of **direct discrimination could capture some cases of algorithmic discrimination where a whole group is consistently impacted**, no matter the criterion used for decision-making.¹⁵³ Going further, **fitting the discriminatory effects of algorithmic bias within one or the other notion raises crucial normative questions** about key concepts of non-discrimination law.¹⁵⁴ In this sense, CAHAL recognised in its 2020 Feasibility Study that **“[t]he increased prominence of proxy discrimination in the context of machine learning may raise interpretive questions about the distinction between direct and indirect discrimination or, indeed, the adequacy of this distinction as it is traditionally understood”.**¹⁵⁵ For example what can be considered a “neutral” criterion for decision-making in light of existing feedback loops and redundant encoding issues? Is algorithmic discrimination, which feeds structural inequality into individual decision-making, a collective or individual form of unfair treatment? Should the user of an algorithm be considered a perpetrator when a machine autonomously “learns” to discriminate? Answers to these questions

151. For example, part-time work is a matter of gender equality where most part-time workers are women. See Tobler, Limits and potential of the concept of indirect discrimination, 24 and Janneke Gerards, ‘Discrimination grounds’, in: Dagmar Schiek, Lisa Waddington and Mark Bell (eds), *Cases, Materials and Text on National, Supranational and International Non-Discrimination Law*, Oxford and Portland, Oregon: Hart Publishing 2007, 33-184.

152. Gerards J and Xenidis R, *Algorithmic discrimination in Europe: Challenges and Opportunities for EU Gender Equality and Non-Discrimination Law* (European network of legal experts in gender equality and non-discrimination / European Commission, 2021).

153. See Adams-Prassl, Binns and Kelly-Lyth, “Directly discriminatory algorithms”, *Modern Law Review* (forthcoming).

154. Gerards J and Xenidis R, *Algorithmic discrimination in Europe: Challenges and Opportunities for EU Gender Equality and Non-Discrimination Law* (European network of legal experts in gender equality and non-discrimination / European Commission, 2021).

155. CAHAL, “Feasibility Study on legal framework on AI design, development and application based on CoE standards” (2020), [13], p. 5.

will determine, in theory, whether the notion of direct or indirect discrimination can be used to capture algorithmic discrimination.¹⁵⁶

Procedural issues: proof, proportionality, responsibility and liability

In practice, however, the opacity of algorithmic decision-making systems means that the evidence necessary to characterise direct discrimination will often be lacking. The information might only become available *ex-post* and might remain partial, so that one might only be able to observe the effects of an algorithmic system after it has been used. For example, if a credit scoring algorithm systematically denies credit to people living with a disability, one might not have access to the criteria used for such a decision but might only be able to observe a pattern of rejection in relation to applicants with a disability. Similarly, one might not be able to access information regarding the entire pool of applicants, so that there might not be any certainty regarding potential applicants with a disability who have been granted credit or other applicants who received a rejection.

Proof issues: For potential applicants, the opacity of algorithmic decisions amounts to substantial barriers to redressing discrimination. Information asymmetries between users and subjects of algorithmic decision-making or decision-support systems mean that isolated end users will not have the capacity to monitor the impact of algorithmic decisions on groups of other end users. They will not be able to access information about the decision-making criteria either. Even in potential cases of indirect algorithmic discrimination, the absence of transparent and meaningful information on relevant decision criteria and victims' lack of a birds-eye view of decisions taken could prevent awareness that discrimination has occurred. This can eventually preclude any legal action from even being started.

Existing rules on the burden of proof are meant to support applicants when bringing cases to court: once a *prima facie* case of discrimination has been established by the applicant, in principle the burden of proof shifts to the defendant, who is responsible for showing that the difference in treatment is justified. **Yet, legal issues still arise: How to provide enough elements, and which type of information to adduce, to make a *prima facie* case of discrimination so as to trigger the shift of the burden of proof onto**

156. See Gerards J and Xenidis R, Algorithmic discrimination in Europe: Challenges and Opportunities for EU Gender Equality and Non-Discrimination Law (European network of legal experts in gender equality and non-discrimination / European Commission, 2021) and Xenidis R, 'Tuning EU Equality Law to Algorithmic Discrimination: Three Pathways to Resilience' (2021) 27 Maastricht Journal of European and Comparative Law 736.

the defendant? In the algorithmic context, information asymmetries might defeat even the possibility to show discrimination *prima facie*.¹⁵⁷

Proportionality test: Once a differential treatment between similarly situated persons or the absence thereof between differently situated persons has been established, judges must conduct a proportionality test to assess whether it can be objectively justified. This two-step test aims to find whether the practice fulfils a legitimate aim, and whether the means employed are reasonably proportionate to the aim pursued.¹⁵⁸ Answering these questions lead to considerable legal uncertainty because of the necessity for judges to assess technical trade-offs that might not be accessible or intelligible to them (e.g., which fairness metrics were to be used? How to balance trade-offs between various definitions of equity?¹⁵⁹ How to balance accuracy vs fairness? etc.).¹⁶⁰ The technical barriers arising here could contribute to shielding algorithmic decision-making systems from judicial review. In these conditions, recent research points towards a permissive application of the proportionality test in the context of algorithmic opacity.¹⁶¹

Responsibility and liability: The question of responsibility and liability for algorithmic discrimination is thorny. Some commentators argue that the law should allow for “an extension of the grounds for defence of respondents [which] could allow them to establish that biases were autonomously developed by an algorithm”.¹⁶² However, such an argument raises the difficult

157. In this context, the Council of Europe Recommendation CM/Rec(2020)1 of the Committee of Ministers to member States on the human rights impacts of algorithmic systems could ease applicants’ difficulty to adduce evidence as it proposes a principle of contestability: ‘As a necessary precondition, the existence, process, rationale, reasoning and possible outcome of algorithmic systems at individual and collective levels should be explained and clarified in a timely, impartial, easily-readable and accessible manner to individuals whose rights or legitimate interests may be affected, as well as to relevant public authorities’.

158. Registry of the European Court of Human Rights, Guide on Article 14 of the European Convention on Human Rights and on Article 1 of Protocol No. 12 to the Convention (30 April 2022) available at: https://www.echr.coe.int/Documents/Guide_Art_14_Art_1_Protocol_12_ENG.pdf (last accessed 22 July 2022).

159. Equity is a philosophical and statistical term used to describe whether an algorithmic system treats different groups fairly. There are different definitions of equity (e.g., all groups get similar rates of false positives and negatives vs the performance of an algorithm is calibrated to be similar for all groups) that can be incompatible with each other. There is no neat overlap between the statistical term ‘equity’ and the legal term ‘equal treatment’.

160. See Binns R, ‘Algorithmic Decision-making: A Guide For Lawyers’ (2020) 25 *Judicial Review* 2.

161. Pablo Martínez-Ramil, ‘Discriminatory algorithms. A proportionate means of achieving a legitimate aim?’ (2022) *Journal of Ethics and Legal Technologies* 4(1).

162. Grozdanovski L, ‘In search of effectiveness and fairness in proving algorithmic discrimination in EU law’ (2021) 58 *Common Market Law Review*, 99.

question of who should be held liable for algorithmic discrimination in the absence of legal personhood of AI systems. Moreover, the distribution of liability between AI providers and users (those deploying them) is another difficulty as both could bear responsibility for a discriminatory system. In light of the many sources of algorithmic bias, from data to model features and implementation, it is nearly impossible to identify a single and precise cause of algorithmic discrimination.

Issues relating to the personal scope of non-discrimination law: the mismatch between algorithmic systems and protected grounds of discrimination

The last set of challenges that arises concerns the lack of overlap between the personal scope of non-discrimination legal provisions and the idiosyncratic forms of algorithmic subjectivity.

Proxy discrimination and the indirect discrimination route: Research shows that algorithmic discrimination takes place even when protected characteristics are removed from a given dataset. This is because algorithmic profiling relies on data points which, combined, can lead to clustering that overlaps with protected groups. For instance, commuting time between home and workplace or postcode could lead to inferences about socio-economic status and ethnicity given the existing spatialization of socio-economic and racial inequalities.¹⁶³ In particular, **redundant encoding** issues arise when variables in a dataset correlate with a protected category, for instance commuting time and ethnic background, which can be inferred by machine learning algorithms. This combines with issues of **feedback loops**, which describe situations where a system relies on data arising from past discrimination as a basis for predictions. Algorithmic discrimination is therefore very likely to take the form of **proxy discrimination**.

Article 14 ECHR bans discrimination “on any ground such as sex, race, colour, language, religion, political or other opinion, national or social origin, association with a national minority, property, birth or other status”. Proxy discrimination, for example based on **behavioural data such as screen time, wifi usage, geolocalisation data**, etc., could therefore be captured under Art. 14 ECHR **via the indirect discrimination route**, by showing a strong

163. See Williams BA, Brooks CF and Shmargad Y, 'How Algorithms Discriminate Based on Data They Lack: Challenges, Solutions, and Policy Implications' (2018) 8 Journal of Information Policy 78.

disadvantageous effect based on one of the grounds explicitly listed.¹⁶⁴ **The problem is that such proxy discrimination might escape the legal protection against discrimination because of the procedural difficulties** exposed in the above section.¹⁶⁵

“New” algorithmic groups and the notion of “other status”: In addition, algorithms can generate new categorizations based on seemingly innocuous characteristics, such as web browser preferences or apartment number, or more complicated categories combining many data points. For example, an online store may find that most consumers using a certain web browser pay less attention to prices; the store can charge those consumers extra. Despite not corresponding to criteria protected under non-discrimination law, some of these algorithmic groups might deserve legal protection, for example if patterns of algorithmic differentiation expose them to systematic socio-economic disadvantage.

Where discrimination against algorithmic groups does not overlap with categories explicitly protected by Art. 14 ECHR, **the open-ended list of protected grounds in Art. 14 and the flexible approach of the European Court of Human Rights (the Court) towards protecting “new grounds” arguably provides an avenue for protection.**¹⁶⁶ It has been argued that “semi-open” anti-discrimination clauses such as Art. 14 ECHR provide better solutions for redressing algorithmic discrimination than fully closed discrimination provisions such as in EU secondary law.¹⁶⁷ For instance, the Court has protected groups on the basis of their professional status or place of residence.¹⁶⁸ This open-ended approach, based on the notion of **“other status”**, could facilitate extending the coverage of new algorithmic groups under Art. 14 ECHR. Yet,

164. Proxy discrimination could in certain cases be treated as direct discrimination, depending on how the scope and boundaries of protected groups are delineated. For a discussion of this problem within the notion of direct discrimination in the EU context, see Xenidis R, ‘Tuning EU Equality Law to Algorithmic Discrimination: Three Pathways to Resilience’ (2021) 27 Maastricht Journal of European and Comparative Law 736.

165. See e.g., Anton Vedder & Laurens Naudts (2017) Accountability for the use of algorithms in a big data environment, *International Review of Law, Computers & Technology*, 31:2, 206-224 and Naudts, L. (2019). How Machine Learning Generates Unfair Inequalities and How Data Protection Instruments May Help in Mitigating Them. In R. Leenes, R. van Brakel, S. Gutwirth & P. De Hert (Editors), *Data Protection and Privacy: The Internet of Bodies* (Computers, Privacy and Data Protection).

166. See Gerards, Janneke, and Frederik Zuiderveen Borgesius. “Protected Grounds and the System of Non-Discrimination Law in the Context of Algorithmic Decision-Making and Artificial Intelligence.” *Colorado Technology Law Journal*, forthcoming (2020).

167. *Ibid.*

168. See *Van der Musselle v. Belgium* Application no. 8919/80 (European Court of Human Rights, 23 November 1983) and *Carson and Others v United Kingdom* (2010), [70]-[71].

this poses the question of the **normative limits of anti-discrimination law**: what are the contours of its mandate? What kinds of injustices is it meant to address?

Furthermore, **some algorithmic clusters lack social salience and are therefore difficult to depict as groups deserving protection from discrimination law**.¹⁶⁹ The “new” algorithmic groups emerging from intangible algorithmic clustering are subject to distinctions that have very tangible socio-economic effects and could consolidate into “**emergent**” **structural discrimination**.¹⁷⁰ By contrast to socially salient algorithmic groups, such distinctions will systematically escape ECHR equality law. Recent scholarship has proposed extending the scope of anti-discrimination law to cover such harmful algorithmic distinctions.¹⁷¹

A last problem pertaining to the personal scope of ECHR equality law arises when **algorithmic decision-making blurs the lines between the individual and the group**. In particular, group-based patterns are used to make decisions about individuals. This presupposes that membership into given algorithmic groups might be ascribed to individuals even when this is not factually correct. For instance, if a user displays the typical web traffic patterns of a woman between 25 and 30 residing in an urban environment, that gender and age identity might be assigned to them to serve as a basis for further decision-making. If that ascribed algorithmic cluster does not match the real user’s identity, the user will not have any opportunity to correct the results of algorithmic profiling and ensuing treatment. However, if that user experienced gender-based discrimination, for example higher health insurance prices, they could claim “**discrimination by association**”, a notion recognised by the Court in 2008.¹⁷²

Intersectional discrimination: Finally, algorithmic discrimination is likely to be intersectional in nature, that is to involve several discrimination grounds

169. See Matthias Leese, The new profiling: Algorithms, black boxes, and the failure of anti-discriminatory safeguards in the European Union, 45 SECURITY DIALOGUE 494–511, 501 (2014); Monique Mann & Tobias Matzner, Challenging algorithmic profiling: The limits of data protection and anti-discrimination in responding to emergent discrimination, 6 BIG DATA & SOCIETY, 5–6 (2019).

170. Ibid.

171. See Wachter S, ‘The Theory of Artificial Immutability: Protecting Algorithmic Groups Under Anti-Discrimination Law’ (2022) Tulane Law Review (forthcoming).

172. Molla Sali v. Greece Application no. 20452/14 (European Court of Human Rights, 19 December 2018), [141].

or vectors of disadvantage.¹⁷³ Because of the granularity of algorithmic profiling, AI systems are able to infer several protected social memberships and potentially **cluster users according to different problematic classifications**. For example, algorithmic profiles might contain information regarding gender, age, ethnic background, religious beliefs, sexual orientation or gender identity based on the analysis of online behaviours, consumer preferences, etc. Identifying and redressing intersectional cases of algorithmic discrimination proves even more challenging than single-axis cases because of the lack of disaggregated equality data, which does not allow comparing potential disparities between algorithmic outputs and the actual situation of intersectionally marginalised groups.¹⁷⁴ Debiasing approaches also show limits when it comes to redressing the discriminatory consequences of biases affecting intersectional minorities.¹⁷⁵ Against this background, intersectional discrimination has often fallen through the cracks of judicial redress. Although the Court has successfully (albeit implicitly) grappled with intersectional discrimination in a case like *BS v Spain*,¹⁷⁶ it has failed to recognise it explicitly and to redress it in others like *SAS v France* or *Garib v The Netherlands*.¹⁷⁷ This **lack of robust legal framework against intersectional discrimination**, often due to formalistic comparison-based conceptions of equality, will prove particularly problematic in the context of algorithmic discrimination.

II. Privacy and data protection law: Fairness and accuracy

In addition to legal instruments pertaining to equality and discrimination, **privacy and data protection law can also be leveraged to tackle algorithmic discrimination**. The concept of fairness in privacy law relates to an

173. The explanatory memorandum to ECRI's General Policy Recommendation No. 14, [1] defines intersectional discrimination as "a situation where several grounds interact with each other at the same time in such a way that they become inseparable and their combination creates a new ground". See also Gerards J and Xenidis R, *Algorithmic discrimination in Europe: Challenges and Opportunities for EU Gender Equality and Non-Discrimination Law* (European network of legal experts in gender equality and non-discrimination / European Commission, 2021).

174. Data categorisation might also be problematic and lack representativeness, with consequences on attempts to fix algorithmic discrimination. See Ruberg, B. and Ruelos, S., 'Data for queer lives: How LGBTQ gender and sexuality identities challenge norms of demographics' (2020) *Big Data & Society*.

175. Balayn A and Gürses S, *Beyond Debiasing: Regulating AI and its inequalities* (European Digital Rights 2021), 62-63.

176. *B.S. v. Spain* Application no. 47159/08 (European Court of Human Rights, 24 July 2012).

177. See e.g., *S.A.S. v. France* Application no. 43835/11 (European Court of Human Rights, 1 July 2014) or *Garib v. The Netherlands* Application no. 43494/09 (European Court of Human Rights, Grand Chamber, 6 November 2017).

organisation's intent to use personal information in good faith, with the intention of balancing the interests of data controllers and data subjects (the individuals). There is general agreement for example that the processing of personal information which is beyond an individual's knowledge/expectation would lead to an unfair situation in the eyes of privacy regulators. **However, the idea of fairness can have many possible nuances: non-discrimination, fair balancing, procedural fairness, bona fide, etc.**

The relation between discrimination and (un)fairness can be found in many legislative acts, proposals, and policy documents across the globe. Convention 108+, alongside the General Data Protection Regulation (GDPR) and many other privacy laws around the world, states that, in order to ensure fair and transparent processing in respect of the data subject, the controller should use appropriate mathematical or statistical procedures for profiling, and implement technical and organisational measures appropriate to prevent potential risks for the interests and rights of the data subject. Risks may include discrimination on the grounds of racial or ethnic origin, political opinion, trade union membership, genetic status or sexual orientation.

Fairness is an overarching principle which requires that personal data shall not be processed in a way that is detrimental, discriminatory, unexpected or misleading to the data subject. It can be argued that fairness in privacy law relates to the **need to address the power imbalance between data subjects (individuals) and the digital ecosystem** and, for this reason, in recent times, privacy law has been leveraged quite extensively to deal with the harms of AI and algorithmic decision making, as outlined in a report issued by the Future Privacy Forum.¹⁷⁸ The report highlights actions taken by Data Protection Authorities including detailed transparency obligations about the parameters that led to an individual automated decision, a broad reading of the fairness principle to avoid situations of discrimination, and strict conditions for valid consent in cases of profiling and automated decision making.

For the purpose of this study, we are looking into two elements of fairness from a privacy standpoint:

- ▶ **Fairness as procedures:** transparency and fairness are inextricably linked because it is arguable that opening the source code to external scrutiny or providing a meaningful explanation on the processing of personal information by the AI system could lead to identification of bias and its root causes, and thus a positive increase in public accountability. For example, The Italian *Corte di Cassazione* issued a sentence in 2021

178. AUTOMATED DECISION-MAKING UNDER THE GDPR – A COMPREHENSIVE CASE-LAW ANALYSIS, Future Privacy Forum, available at: <https://fpf.org/blog/fpf-report-automated-decision-making-under-the-gdpr-a-comprehensive-case-law-analysis/>

stating that a data subject's consent cannot be deemed valid if the algorithm is not transparent as the data subject would not be able to understand what they are consenting to.¹⁷⁹ This case was welcomed by the Italian privacy regulator, Garante, as a demonstration of how the privacy law (and the GDPR in this case) is fit for upholding individuals' rights in the age of AI.

- **Fairness as the protection of individual vulnerabilities:** in privacy law, fairness is often conceived as a corrective tool for rebalancing asymmetric or unbalanced relationships between organisations and individuals. Take for example the case of algorithmic platforms where the French Conseil d'Etat (as rephrased by Commission Nationale de l'Informatique et des Libertés) affirms that "fairness consists of ensuring, in good faith, the search engine optimisation (SEO) or ranking service, without seeking to alter or manipulate it for purposes that are not in the users' interest".¹⁸⁰ On a more general level, in the algorithmic environment, "fairness could well represent a solution to the problem of *unbalanced relations* between controllers of algorithms and users".¹⁸¹

For many countries both within and outside of Europe, the modernization of Convention 108 – with the introduction of **new rights for data subjects in algorithmic decision-making contexts**, particularly in connection with artificial intelligence – represents a common ground, as the treaty serves as a borderline standard for how countries should go about protecting the privacy rights of their citizens in the age of AI. The GDPR, which has many similarities with **Convention 108 +** (although the Council of Europe has a much wider reach and territoriality than the EU) also contains provisions to support individual rights in the context of AI and algorithms, including the renowned Article 22, which safeguards individuals from automated decision-making.

There are several other safeguards that apply to such data processing activities, notably those stemming from the general data processing principles in Article 5, the legal grounds for processing in Article 6, the rules on processing special categories of data (such as biometric data) under Article 9, specific transparency and access requirements regarding algorithmic decision-making (ADM) under Articles 13 to 15, and the duty to carry out data protection impact assessments in certain cases under Article 35.

179. Corte di Cassazione, Civile Ord. Sez. 1 Num. 14381, ItalggiureWeb, 25 May 2021 available at: <http://www.italgiure.giustizia.it/xway/application/nif/clean/hc.dll?verbo=attach&db=snciv&id=../20210525/snciv@s10@a2021@n14381@tO.clean.pdf> (last accessed: 26 May 2021)

180. Conseil d'État, "Le Numérique et les droits fondamentaux", 2014, pp. 273 and 278-281.

181. Understanding algorithmic decision-making: Opportunities and challenges, available at: [https://www.europarl.europa.eu/RegData/etudes/STUD/2019/624261/EPRS_STU\(2019\)624261_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/STUD/2019/624261/EPRS_STU(2019)624261_EN.pdf)

However, there are **limitations** in current privacy instruments when it comes to AI and algorithmic decision making, including:

- ▶ Exercising data subjects' rights in the context of AI and algorithmic decision making is rather complex. For example, even with the guidance of Data Protection Working Party 29 on automated individual decision-making and profiling, the assertion of GDPR Article 22 ("**solely automated**, and **legal or similarly significant effects**") presents practical challenges.
- ▶ Transparency of algorithmic management is the first step towards genuine accountability. However, **transparency and explainability** requirements in relation to bias mitigation raise questions around the **intersection of privacy and trade secret laws**. Importantly, algorithmic systems need to meet a certain threshold of accessibility and intelligibility, whether by internal or external auditors, a regulator, or a tribunal. However, a company's own algorithm may also be covered by trade secrets legislation. There are interesting developments in this sense thanks to the emergence of Secure Multi Party Computation that may enable an AI to be interrogated without having access to the actual code. Nevertheless, that is still a long way off. Current regulatory efforts go in this directive, for instance the Council of Europe Recommendation CM/Rec(2020)1 on the human rights impacts of algorithmic systems provides that '[t]he legislative frameworks for intellectual property or trade secrets should not preclude such transparency, nor should States or private parties seek to exploit them for this purpose' and that '[c]onfidentiality considerations or trade secrets should not inhibit the implementation of effective human rights impact assessments'.

III. AI sectoral regulations: strengths and limits for promoting equality and addressing discrimination

In addition to discrimination, privacy and data protection laws, sectoral regulations will also be relevant for addressing algorithmic discrimination.

The Council of Europe is currently developing regulation that would address algorithmic discrimination as part of an effort to promote human rights, democracy and the rule of law. This would take the form of a legally binding transversal instrument addressing issues in the public sector as well as binding and non-binding sectoral regulations.¹⁸² In 2020 CAHAI prepared a

¹⁸². See CAHAI, "Feasibility Study on legal framework on AI design, development and application based on CoE standards" (2020), [54].

“Feasibility Study on a legal framework on AI design, development and application based on Council of Europe standards”, which recognises that “AI systems [can] be used in a way that perpetuates or amplifies unjust bias, also based on new discrimination grounds in case of so called ‘proxy discrimination’”.¹⁸³ At the same time, CAHAI considers that “AI systems can foster and strengthen human rights more generally, and contribute to the effective application and enforcement of human rights standards”, for instance “by detecting biased (human or automated) decisions, monitoring representation patterns of different people or groups (for example women in the media) or analysing discriminatory structures in organisations”.¹⁸⁴

In its 2021 document **“Possible elements of a legal framework on artificial intelligence, based on the Council of Europe’s standards on human rights, democracy and the rule of law”**, CAHAI recommends including “a provision on respect of *equal treatment and non-discrimination* of individuals in relation to the development, design, and application of AI systems to avoid unjustified bias being built into AI systems and the use of AI systems leading to discriminatory effects” in the legally binding transversal Framework Convention on AI regulation which is currently under preparation.¹⁸⁵

CAHAI also proposes complementary regulation for the public sector, where it recommends that “documentation and logging processes” pertaining to the development of the system “should be meticulously kept to ensure transparency and traceability”. It recommends that “[a]dequate test and validation processes, as well as data governance mechanisms should be put in place” to assess risks “of unequal access or treatment, various forms of bias and discrimination, as well as the impact on gender equality”.¹⁸⁶

As other sectoral regulations are envisaged in Europe, it is important to flesh out the **added value of regulating AI at Council of Europe level**. Arguably, regulation by the Council of Europe can have **strong influence worldwide** due to the broad membership of the Council of Europe, its distinctive human rights-based approach and the fact that the instrument would be open for ratification to non-state parties as well. The CAHAI’s “Possible Elements” document point towards minimum standards and an approach focused on the public sector, in line with the European Convention on Human Rights mechanism, which differs from the “market approach” taken by the EU in its draft

183. Committee on Artificial Intelligence, “Possible elements of a legal framework on artificial intelligence, based on the Council of Europe’s standards on human rights, democracy and the rule of law”, Council of Europe (2022), [13]

184. *Ibid.*, [20].

185. *Ibid.*, [27]

186. *Ibid.*, [60]

EU AI Act.¹⁸⁷ A commonality between the two regulations would be the risk-based approach they both adopt to AI systems.¹⁸⁸ Yet the Council of Europe has the potential to foster a distinct **human-rights-based approach to AI and algorithmic technologies**.

Sectoral regulation of AI is also currently underway in the EU. The draft **EU AI Act** follows a risk-based approach and classifies AI systems as **“high-risk”** if they are deployed in the following areas: biometric identification and categorisation of natural persons, management and operation of critical infrastructure (road traffic, water, gas, heating and electricity supply), education and vocational training, employment, workers management and access to self-employment, access to and enjoyment of essential private services and public services and benefits, law enforcement, migration, asylum and border control management, administration of justice and democratic processes. AI systems that present an **“unacceptable risk”**, are prohibited for example “practices that have a significant potential to manipulate persons through subliminal techniques beyond their consciousness or exploit vulnerabilities of specific vulnerable groups such as children or persons with disabilities in order to materially distort their behaviour in a manner that is likely to cause them or another person psychological or physical harm”. AI systems that present a **limited risk** are subjected to specific transparency obligations and those with **low or minimal risk** to codes of conduct.

Although the EU AI Act foresees promising transparency obligations with a view to bias mitigation, in particular in relation to training data and decision criteria,¹⁸⁹ several **criticisms** have been put forward regarding the way in which the EU AI Act proposes to ensure that fundamental rights are respected. For example, it approaches AI systems from a product liability perspective and thus **does not foresee complaint mechanisms** that would enable **victims of algorithmic discrimination or NGOs with a legitimate interest to request that changes are made to these systems after their deployment** in compliance with anti-discrimination law.¹⁹⁰ Moreover, commentators have criticised

187. See Marten Breuer, “The Council of Europe as an AI Standard Setter” *Verfassungsblog* (4 April 2022) available at: <https://verfassungsblog.de/the-council-of-europe-as-an-ai-standard-setter/>.

188. See Committee on Artificial Intelligence, “Possible elements of a legal framework on artificial intelligence, based on the Council of Europe’s standards on human rights, democracy and the rule of law”, Council of Europe (2022), [19].

189. See in particular Art. 10 on Data and data governance of the EU AI Act.

190. See Joan Lopez Solano, Aaron Martin, Siddharth de Souza and Linnet Taylor, “Governing data and artificial intelligence for all Models for sustainable and just data governance” (Panel for the Future of Science and Technology, European Parliamentary Research Service 2022), 52 available at: [https://www.europarl.europa.eu/RegData/etudes/STUD/2022/729533/EPRS_STU\(2022\)729533_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/STUD/2022/729533/EPRS_STU(2022)729533_EN.pdf).

the **absence of legal obligations** for providers and users of AI systems to **conduct *ex-ante* human rights impact assessments**.¹⁹¹ The absence of any equality mainstreaming clause or positive obligation requiring AI and algorithmic systems to promote equality is also regrettable. **These are aspects on which the Council of Europe instrument should focus in order to create complementarity with the EU AI sectoral regulations and to ensure that its human rights mandate is at the core of the new legal provisions.**

In 2022, the European Commission **proposed a new directive on adapting non-contractual civil liability rules to artificial intelligence**.¹⁹² The proposal aims to ‘enable effective private enforcement of fundamental rights and preserve the right to an effective remedy where AI-specific risks have materialised’, including non-discrimination. The Commission explains that its proposal ‘complements other strands in the Commission’s AI policy based on preventive regulatory and supervisory requirements aimed directly at avoiding fundamental rights breaches (such as discrimination)’. While it ‘does not create or harmonise the duties of care or the liability of various entities whose activity is regulated under [non-discrimination law] and, therefore, does not create new liability claims or affect the exemptions from liability under [non-discrimination law]’, it ‘introduces **alleviations of the burden of proof for the victims of damage caused by AI systems** in claims that can be based on national law or on these other EU [non-discrimination law]’. This Study suggests that **these rules could be used as inspiration by the Council of Europe for facilitating applicants’ access to justice** with regard to claims of algorithmic discrimination, in particular in relation to issues of proof and evidence.

As explained in Section 3 below, future AI sectoral regulations at Council of Europe level should also include a **legal obligation for AI and algorithmic systems to promote equality**. Norwegian equality legislation could offer a useful yardstick in this context as it foresees equality promotion as a legal obligation.¹⁹³

191. See *ibid.*

192. European Commission, Proposal for a Directive of the European Parliament and of the Council on adapting non-contractual civil liability rules to artificial intelligence (AI Liability Directive) COM(2022) 496 final.

193. See Chapter 4 of the Norwegian Act relating to equality and a prohibition against discrimination (Equality and Anti-Discrimination Act), available at: https://lovdata.no/dokument/NLE/lov/2017-06-16-51#KAPITTEL_4.

Section 3

Promoting equality in and through the use of AI: the role of positive action and positive obligations

While the previous section highlighted relevant legal and policy instruments and Council of Europe level and at EU and international level, it has also pointed at problematic gaps, shortcomings and uncertainties in the applicability of these instruments to the problem of algorithmic discrimination. As shown in this section, avenues for fixing these issues should promote a **paradigm shift**. First, it could be recommended that **existing rules should be revisited in light of the new power and information asymmetries inherent in algorithmic technologies**. Second, we recommend that **positive action and positive obligations be used as an avenue for crafting a legal obligation to prevent discrimination and promote equality in and through the use of algorithmic systems**. Taking these two steps would elevate ‘equality by design’ as a prominent feature of the Council of Europe’s human-rights-based approach to algorithmic discrimination,

I. Revisiting existing rules in light of new power asymmetries

This section aims to outline avenues for responding to the issues highlighted in Section 2 in relation to the applicability of existing legal provisions.

First, in light of existing research which has shown that in the absence of safeguards, algorithmic bias systematically pervades algorithmic decisions, a **presumption of algorithmic bias** could be posited where no preventive measures have been taken by users of algorithmic systems. This is justified by the pervasiveness of bias in the design process of AI systems, ranging from biases in data collection and datasets to biases in problem design,

algorithmic models and implementation of AI recommendations.¹⁹⁴ As argued by Eubanks, “when automated decision-making tools are not built to explicitly dismantle structural inequalities, their increased speed and vast scale intensify them dramatically”.¹⁹⁵ In other terms, algorithmic discrimination is very likely to arise where no safeguards have been put in place. When it perpetuates inequality, the use of biased AI systems should be equated with actively enacting structural disadvantage and amplifying the unfair distribution of valuable social goods. The **foreseeability of discriminatory harms arising from algorithmic bias** thus justifies conceptualising algorithmic discrimination as a form of **negligence**. Drawing from Moreau’s work on discrimination and tort-based theories of discrimination law,¹⁹⁶ it is possible to derive a **social responsibility for users of algorithmic systems to take reasonable action to prevent the aggravation of discrimination** in society. This approach resonates with the discussions currently taking place in the EU context and in particular the Commission’s proposal for a “rebuttable presumption for AI-related damages”.¹⁹⁷

Second, the pervasive use of AI systems establishes **new power and information asymmetries**. It becomes very **difficult for subjects of algorithmic decisions to identify discrimination** due to a combination of personalisation, automation and opacity of decision-making processes. Comparison with similarly placed individuals and social interactions are important heuristic devices when it comes to acquiring presumptions of discrimination. Yet, reading social cues or comparing oneself to other loan applicants in the context of an online credit service becomes impossible.¹⁹⁸ This information asymmetry makes it difficult to suspect discrimination in the first place. Even when suspicion arises, **collecting evidence is a further challenge** because

194. Grozdanovski suggests that it is possible to read the existence of such a presumption in the EU White paper on Artificial Intelligence, see Grozdanovski L, ‘In search of effectiveness and fairness in proving algorithmic discrimination in EU law’ (2021) 58 Common Market Law Review.

195. Eubanks V, *Automating inequality: how high-tech tools profile, police, and punish the poor* (First edition, edn, St. Martin’s Press 2018).

196. See Sophia Moreau, ‘Discrimination as negligence’ (2010) 40 *Canadian Journal of Philosophy* 123; Oppenheimer DB, ‘Negligent Discrimination’ (1993) 141 *University of Pennsylvania law review* 899.

197. See in this sense Luca Bertuzzi, “LEAK: Commission to propose rebuttable presumption for AI-related damages” (Euractiv, 2022) available at: <https://www.euractiv.com/section/digital/news/leak-commission-to-propose-rebuttable-presumption-for-ai-related-damages/>.

198. In private sector application of AI-powered systems in particular, individuals may not necessarily be informed that an algorithmic system is involved in a decision that pertains to them. Thus, it becomes even more challenging to know when to request human review, or to pay attention to the particularities of potential algorithmic discrimination.

decisions or the algorithmic recommendations supporting them are not readily available to consult and often not disclosed by users of algorithmic decision-making systems. Hence, **presenting proof to establish a presumption of discrimination in courts is a key legal challenge**. Even though the shift of the burden of proof can help mitigate the power asymmetries created by opaque algorithmic systems,¹⁹⁹ the threshold to trigger this shift should reflect end users' position and limited access to *prima facie* evidence.

Bringing together the foreseeability of algorithmic bias and existing information asymmetries reveals how the pervasive deployment of AI systems in decision-making processes **upsets the balance between the position of possible victims of discrimination and that of the providers and users of these systems**. While victims are subjected to more pervasive discrimination which they are contemporaneously less able to identify and prove, profit-makers enjoy increased power thanks to AI systems that enhance economic profits while possibly shielding them from liability for their discriminatory consequences due to the legal obstacles listed above. Hence, **the legal framework needs to be adjusted to reflect and integrate the power shifts and imbalances** that derive from the use of AI systems in a vast array of decisions that open or close life opportunities and therefore intensely affect inequality in society.

Revisiting existing rules on the burden of proof can help restore the effectiveness of non-discrimination law in light of new power and information asymmetries between users and subjects of algorithmic decision-making systems. Positing a presumption of algorithmic bias as suggested above would allow **shifting the burden of proof onto the defendant as soon as no preventive measures have been taken**. Such preventive measures could take the form, for instance, of an impact assessment, an audit or a certification of the algorithmic system used, as exposed in the recommendations section. Failure to take adequate preventive measures could then amount to negligence. This mechanism would support potential victims in adducing accessible *prima facie* evidence with a view to shifting the burden of proof onto users. Such an adaptation of the legal framework would also **mainstream positive action and preventive obligations** against algorithmic bias, as further outlined below.

199. See C-109/88 Handels- og Kontorfunktionærernes Forbund I Danmark v Dansk Arbejdsgiverforening, acting on behalf of Danfoss EU:C:1989:383.

Third, the adaptation of existing rules suggested above should be combined with a public supervisory approach.²⁰⁰ **Empowering equality bodies, discrimination ombudspersons and national human rights institutions to monitor the discriminatory impact of algorithmic decision-making and support systems** should be made a priority. This involves providing these institutions with necessary legal rights and investigative powers (e.g., to access datasets and decision criteria), the right resources, but also with capacity to prevent discrimination by co-operating with users of algorithmic decision-making systems – for instance companies using algorithmic decision-making systems to support recruitment procedures – to collect relevant data on the impact of their decisions, and to assist potential victims in relation to obtaining redress. Monitoring could take the form of **situation testing** where these authorities test the outcomes of a given system by comparing results for different groups. For instance, they could submit test CVs or credit applications from majority and minority groups to try and reveal algorithmic discrimination in contexts where companies use ADM systems. They could also conduct **audits** to detect potential bias if granted access to relevant systems. Such **public enforcement** methods could support victims by mitigating existing obstacles to establishing *prima facie* discrimination.

The monitoring function of equality bodies should be supported by **legal obligations around transparency**. Users of algorithmic systems should be required to **provide meaningful and intelligible information on the criteria used for decision-making**. At the moment, the GDPR does not offer a right to explanation.²⁰¹ In the area of goods and services, consumer protection should also be explored as a tool to request information about algorithmic decisions for consumers who have been potentially discriminated against. This could help address the power asymmetries created by the opacity of algorithmic decision-making systems between the subjects of algorithmic decisions and their authors. For this purpose, the recent Council of Europe Recommendation CM/Rec(2020)1 on the human rights impacts of algorithmic systems mentioned in detail in Section 2 of this Study offers a number of interesting pathways.

Fourth, it is necessary to ensure the reviewability of algorithmic systems in light of non-discrimination obligations. Where applicants, lawyers or judges

200. See Xenidis R and Senden L, 'EU Non-discrimination Law in the Era of Artificial Intelligence: Mapping the Challenges of Algorithmic Discrimination' in Bernitz U and others (eds), *General Principles of EU Law and the EU Digital Order* (Wolters Kluwer 2019).

201. See Wachter, Sandra, Brent Mittelstadt, and Luciano Floridi. "Why a right to explanation of automated decision-making does not exist in the general data protection regulation." *International Data Privacy Law* 7.2 (2017): 76-99.

are presented with technical information concerning a specific system, such information is not likely to be intelligible in terms of the system's discriminatory or non-discriminatory nature. **Technical discussions about the adequacy of given fairness metrics and appropriate thresholds for trade-offs between accuracy and equity are difficult to assess from the perspective of legal obligations arising from anti-discrimination law.** In this context, how to ensure that algorithmic decision-making systems undergo a **proportionality test** that guarantees the effectiveness of non-discrimination law? Here again, several solutions can be envisaged, as further articulated in the recommendations section of this Study. On the one hand, **transparency obligations** weighing on the users of algorithmic decision-making systems could guarantee access to an intelligible account of the technical and fairness choices made by developers and users. On the other, **mainstreaming positive action** could lead to a **positive obligation to prevent algorithmic bias** that would **displace the proportionality assessment from the technical to the legal terrain.** What judges would consider, then, would rather be the appropriateness of the preventive measures taken to ward off bias, rather than technical fairness and equity choices.

Finally, we suggest that **liability, as a judicial construct approximating responsibility, should be allocated strategically so as to facilitate access to justice and remedies** in cases of algorithmic discrimination. In the context of the ECHR and other legal instruments at Council of Europe level, where obligations weigh on public authorities, we suggest that **state parties should hold users of AI systems liable for algorithmic discrimination arising from the deployment of their system.** As explained in the Section on recommendations, **this can be complemented with legal obligations for providers to conduct human rights impact assessments *ex-ante* to prevent discriminatory harms.** This will also allow encourage the **documenting of any preventive measures** taken by the provider so as to **ensure that meaningful information can be provided to the user and end-users** of the system in case of legal proceedings.

The approach proposed here, which revolves around a presumption of algorithmic bias, negligence and prevention, could contribute to legal certainty and the effectiveness of the ECHR anti-discrimination provisions by alleviating victims' burden of proof, fostering preventive safeguards, clarifying the allocation of liability and helping better define available justifications for defendants. All in all, we suggest that **a more substantive approach to equality should drive the interpretation of anti-discrimination provisions** to safeguard their effectiveness in the context algorithmic discrimination.

II. An obligation to promote equality in and through the use of algorithmic systems: the role of positive action and positive obligations

This report has shown how AI systems, without the right guardrails and controls, can lead to further exclusion of women and vulnerable groups. Notwithstanding the discriminatory potential of AI, researchers and developers have explored the opportunities offered by AI for identifying and redressing inequality. This requires a **paradigm shift where baselines for software design and deployment are systematically called into question and checked in relation to their inclusionary or exclusionary impact**. In other terms, the deployment of a new AI system should be “purposeful and intentional in its inclusivity” and “must empower communities and present a benefit to all of society”.²⁰² This requires a set of obligations on companies to require them to do so, and a set of pre-market and post-release controls. Below, we argue that such a paradigm shift requires the vast array of available positive action measures including awareness-raising, promotion-based measures, temporary special measures and quota to be utilised for equality, diversity and inclusion purposes across the board.

Rooting out bias and inequality requires a conscious, arguably political and social choice. In the first instance, it should be recognised that AI systems are not neutral but reproduce and amplify structural inequality and the systems of exclusion and disadvantage that are institutionalised in society. This necessitates stepping away from a perpetrator’s perspective on discrimination and instead acknowledging that majority norms and unquestioned assumptions underlying software development and deployment lead to the needs of women and minority groups not being accommodated.²⁰³ Assuming that a system will equally cater for various groups will *de facto* prevent minority groups from benefitting from AI applications and related opportunities to the same extent as other groups. Therefore, **substantive equality**

202. Renee Cummings, “This is how AI can support diversity, equity and inclusion”, World Economic Forum, available at: <https://www.weforum.org/agenda/2022/03/ai-support-diversity-equity-inclusion/>. See also Equality Now, A Call For An Intersectional Feminist Informed Universal Declaration On Digital Rights, available at: https://www.equalitynow.org/news_and_insights/universal-declaration-on-digital-rights/.

203. For a powerful account of the perpetrator’s perspective on discrimination vs understanding discrimination as a structural phenomenon, see e.g., Freeman AD, ‘Legitimizing Racial Discrimination Through Antidiscrimination Law: A Critical Review of Supreme Court Doctrine’ (1978) 62 Minnesota Law Review. This has been recognised in law through the concept of indirect discrimination.

and anti-discrimination 'by design' should be placed at the centre of the legal regulation of AI development and deployment.

1) What is positive action?

Positive action, also called temporary special measures or positive measures in the European context, is a range of policies that can be adopted with a view to reaching full or *de facto* equality. It builds up on a critique of formal equality or equality of opportunity that denounces these frameworks' blindness towards the different starting positions of different social groups. For example, giving the same job opportunity to a worker with disability and an able-bodied worker might lead to a higher dropout rate in the first case because no accommodation measure has been taken to ensure that the worker living with a disability is actually able to perform their tasks. Instead, anchoring policies in theories of substantive equality dictates the adoption of special accommodation measures that create conditions where historically disadvantaged groups can participate in society and reap the benefits of that participation to the same extent as privileged groups. Concretely, that would mean ensuring that a worker living with a disability can access a safe and adapted physical and psychological working environment, for instance through special equipment, flexible working hours, etc. So-called transformative equality theories point in the same direction but place more conceptual emphasis on transforming the unequal status quo in the long-term, for example through granting specific and temporal advantages to structurally disadvantaged groups. An example of such equality policies is flexible quota schemes whereby, for example, an employer faced with equally qualified male and female candidates in a recruitment process would give preference to the female candidate where women are under-represented in the professional community at stake.

In the context of the Council of Europe, positive action is not a legal obligation but has for example been encouraged by the European Commission against Racism and Intolerance (ECRI) "as an effective tool for achieving a fair and even playing field in society for members of disadvantaged groups".²⁰⁴ The Committee of Ministers Recommendation CM/Rec(2003)3 on balanced participation of women and men in political and public decision-making from 2003 also encourages the Council of Europe member states to ensure that the representation of either women or men in any decision-making body in

204. European Commission against Racism and Intolerance, Seminar with national specialised bodies to combat racism and racial discrimination on positive action: explanatory note (2007), available at: <https://rm.coe.int/seminar-with-national-specialised-bodies-to-combat-racism-and-racial-d/16808b54b0>.

political or public life should not fall below 40%.²⁰⁵ All the Council of Europe member states have also ratified the UN Convention on the Elimination of All Forms of Discrimination against Women, which gives clear support to positive action by stating that “adoption by States Parties of temporary special measures aimed at accelerating de facto equality between men and women shall not be considered discrimination as defined in the present Convention, but shall in no way entail as a consequence the maintenance of unequal or separate standards; these measures shall be discontinued when the objectives of equality of opportunity and treatment have been achieved”.²⁰⁶ In the EU, non-discrimination law allows for special measures in the framework of positive action within certain limits such as the prohibition on strict quota that would give an automatic preference to under-represented groups and the need for special measures to aim to transform the status quo in the long run.²⁰⁷ The definition of positive action in the context of the Council of Europe and the European Convention on Human Rights is similar. The concept of “temporary special measures” is often used. ECRI’s General Policy Recommendation no. 7 for example indicates that “[t]he law should provide that the prohibition of racial discrimination does not prevent the maintenance or adoption of temporary special measures designed either to prevent or compensate for disadvantages suffered by persons [from protected groups] or to facilitate their full participation in all fields of life”.²⁰⁸ It also states that “[t]hese measures should not be continued once the intended objectives have been achieved”.²⁰⁹

2) Positive obligations under the ECHR

To approach the question of how to promote equality in and through the use of AI, the legal basis exposed above, which authorises positive action, can be considered together with another important specific feature of the

205. Council of Europe, CM/Rec(2003)3 on balanced participation of women and men in political and public decision-making adopted by the Committee of Ministers of the Council of Europe, (12 March 2003), available at https://search.coe.int/cm/Pages/result_details.aspx?ObjectID=09000016805e0848

206. United Nations General Assembly, Convention on the Elimination of All Forms of Discrimination against Women, article 4, paragraph 1 (18 December 1979) available at <https://www.ohchr.org/en/instruments-mechanisms/instruments/convention-elimination-all-forms-discrimination-against-women>

207. For a detailed account, see Raphaële Xenidis and Hélène Masse-Dessen, ‘Positive action in practice: some dos and don’ts in the field of EU gender equality law’ (2018) 2 European equality law review 36.

208. ECRI General Policy Recommendation No. 7 on National Legislation to Combat Racism and Racial Discrimination (2002), [5].

209. Ibid.

ECHR, namely the notion of **positive obligations**. Positive obligations entail that states have, in certain circumstances, the duty to actively take measures to achieve equality and prevent discrimination.²¹⁰ This goes further than limited passive or negative duties not to discriminate because it implies taking preventive action against discrimination or positive action measures to promote equality as a means to comply with Article 14 ECHR.

In its General Policy Recommendations No. 7, ECRI specifically endorses positive obligations to promote equality and prevent discrimination in the form of constitutional provisions, duties for public authorities, as well as obligations for public bodies to condition “the awarding of contracts, loans, grants or other benefits” to the respect of the positive obligation to promote equality and prevent discrimination.²¹¹ This can be used as a legal basis to create an equality mainstreaming obligation in the context of AI use by public authorities.

Positive obligations and positive action provide an interesting legal basis for utilising AI to promote equality in two regards. On the one hand, it can be argued that positive obligations to prevent discrimination require states to use positive action in order to create safeguards to prevent unlawful algorithmic bias from emerging at any level of the AI lifecycle. On the other hand, positive obligations to promote equality could be interpreted as a requirement for states to invest in using the new opportunities created by AI to better serve disadvantaged communities so that they can fully enjoy the rights guaranteed by the ECHR. The next paragraphs lay out strategies for doing so.

3) Centring positive action

A *sine qua non* condition for using AI for good is positive action. Positive action can take many forms ranging from support measures such as information dissemination among targeted communities, dedicated training and funding programmes, to temporary special measures and flexible quota systems.²¹² For example, key priorities should include **diversifying** educational and professional communities involved with all phases of the development

210. See European Court of Human Rights, Guide on Article 14 of the Convention (prohibition of discrimination) and on Article 1 of Protocol No. 12 (general prohibition of discrimination) (2022), [42-43] available at: https://www.echr.coe.int/Documents/Guide_Art_14_Art_1_Protocol_12_ENG.pdf. See also e.g., European Court of Human Rights, Application no. 34369/97 *Thlimmenos v. Greece* (2 April 2000) and European Court of Human Rights, Application no. 11146/11 *Horváth and Kiss v. Hungary* (29 January 2013).

211. *Ibid.*, [2], [8] and [9].

212. See Christopher McCrudden, Resurrecting positive action (2020) 18(2) *International Journal of Constitutional Law*, 429.

and deployment of AI applications through financial support and awareness-raising efforts. This can be part of a broader effort to attract and retain more women and girls and people from marginalised communities to STEM fields.

Where necessary, temporary special measures and flexible quota schemes should be used to ensure parity and inclusion in educational and professional communities. Positive action measures in the form of e.g., special accommodation and anti-stereotyping measures should aim to render these environments more inclusive so as to retain minority groups in the long-term and reduce drop-out rates. In the same vein, provisions on mainstreaming non-discrimination obligations, including gender mainstreaming, can be regarded as a fruitful legal basis to realise 'equality by design' in the field of AI.

Training should be provided to these communities via a transformation of educational curricula, with ethical issues, legal requirements and social science approaches to discrimination and inequality being part and parcel of higher and professional education. Complementary training should also be provided regularly to experts, stakeholders and professional communities in the AI industry on an *ad hoc* basis or as continuous education. Such training should address structural inequality, gender mainstreaming, and stereotyping. Training should also target other relevant target groups including monitoring bodies (including equality bodies, national human rights institutions, ombudswomen, etc.) and CSOs, legal professionals, and judges dealing with digital rights and discrimination.

An approach centred on substantive equality and positive action might also require adapting existing legal arrangements. Indeed, as the emergence of new technologies shifts power dynamics between users and subjects of AI systems, the justice arrangements and normative dispositions underpinning legal rules become unsettled. Re-balancing such power asymmetries therefore entails adapting the legal architecture. As explained below, rules around the shift of the burden of proof might be eased for victims of algorithmic discrimination via the positing of a presumption of algorithmic bias.²¹³ Such a presumption could arise where users of

213. Not to be confounded with a presumption of algorithmic discrimination because such bias might or might not be discriminatory. For other suggestions on easing the burden of proof in relation to algorithmic discrimination, see Janneke Gerards and Raphaële Xenidis, *Algorithmic discrimination in Europe: Challenges and Opportunities for EU Gender Equality and Non-Discrimination Law* (European network of legal experts in gender equality and non-discrimination / European Commission, 2021) and *AlgorithmAudit*, White Paper: Reversing the burden of proof in the context of (semi-)automated decision-making (2022) available at: <https://drive.google.com/file/d/1RHdqoGVgwwv-FTv8qC9fAlsVl8eUTcR7s/preview>.

an AI system have not put antidiscrimination safeguards in place, i.e. where they have assumed AI systems to be neutral towards protected groups. As described below, valid safeguards could take several forms such as audits, certifications, equality impact assessments. Further details on this proposed legal adaptation are provided in section 4.

4) Using data analytics to detect discrimination

A second possibility for AI to be used for promoting equality is through deploying the power of data analytics to detect discriminatory patterns in the allocation of resources, the dissemination of information, the representation of groups or the performance of given systems. Several examples show that data analytics can also be utilised to unpack bad models and end practices that replicate bias. For instance, AI image recognition technologies could be used to analyse large amounts of data and assess representations of women and minorities across different media sectors ranging from TV programmes to movies, online and physical advertising, etc. In content moderation, AI has been used to detect hate speech in order to report and remove offensive content.²¹⁴ At the same time, it is crucial to prevent that such deployment of AI silences discriminated or minority groups.²¹⁵ Detecting discriminatory language in job ads automatically could also be a way to put AI to the service of the promotion of equality. Going even further, recommender systems could be used to recommend alternative inclusive language to substitute discriminatory content in job ads.

5) AI as a means to serve discriminated people and underserved communities and improve accessibility

Beyond detection, AI systems can also be purposively developed to serve discriminated people, marginalised, at-risk or underserved communities. For instance, AI can be used to improve accessibility to information or existing goods and services. Training automated translation systems on regional or minority languages that are spoken only by a small number of persons would improve access to key services. AI could also serve the promotion of equality in the criminal and policing sector, for instance when put to use to prevent risks of gender-based violence as in Spain with the VioGen software. In the health sector, AI could be used to enhance access to healthcare in

214. European Commission against Racism and Intolerance, General Policy Recommendation No. 15 On Combating Hate Speech CRI(2016)15, [140] available at: <https://rm.coe.int/ecri-general-policy-recommendation-no-15-on-combating-hate-speech/16808b5b01>.

215. See for example the sexist and racist effects of automate content moderation: Gerrard Y and Thornham H, 'Content moderation: Social media's sexist assemblages' (2020) 22 1266.

disenfranchised areas and to improve diagnosing capacities for traditionally under-represented groups.

The condition for such positive usages of AI is however to invest resources into diversifying and training the professional communities involved in developing and using AI and to take positive action measures to ensure that these systems serve marginalised groups. At the same time, technosolutionism should be avoided and AI should not be perceived as a panacea to solve discrimination. It is crucial to remember that social issues require a social approach – not a purely technological one. While AI can certainly be developed and used for the promotion of equality, including gender equality, it is important to view it as a complementary tool in the framework of well-funded and carefully thought-through equality policies. This requires a conscious shift of approach.

Section 4

Recommendations

Policy recommendations: Towards human-rights-based approach to AI

In light of the legal gaps highlighted in this Study and the complexity of AI-driven and algorithmic discrimination, addressing the problem requires a multi-faceted human rights-based approach. In addition to its work to develop a general convention on artificial intelligence, the Council of Europe should aim to be a leading standard-setter in the specific field of equality through the preparation of a more specific Committee of Ministers Recommendation on AI, equality, including gender equality, and discrimination, drafted by an expert committee under the Gender Equality Commission (GEC) and the Steering Committee on Anti-Discrimination, Diversity and Inclusion (CDADI).

Four complementary avenues for regulatory and policy intervention are identified below. The recommendations of this Study should be read in line with provisions of the general convention under preparation with which they are complementary. They build on, and are fully in line with, Recommendation CM/Rec(2020)1 on the human rights impacts of algorithmic systems adopted in 2020 by the Committee of Ministers.²¹⁶

Actions are primarily addressed to member states for deployment in the public sector – for instance in relation to access to justice, legal redress, democratic participation, public awareness-raising and capacity-building – but many recommendations should apply to the private sector as well, including prevention, transparency and accountability measures and suggestions aimed at improving diversity, inclusion, representation and participation. In addition, the recommendations below do not focus on any specific bias source, but rather address the risks of discrimination in algorithmic usages. Therefore, they pertain to the entire lifecycle of AI systems, from design and modelling to training (including quality and representativeness in data

216. Recommendation CM/Rec(2020)1 of the Committee of Ministers to member States on the human rights impacts of algorithmic systems (Adopted by the Committee of Ministers on 8 April 2020 at the 1373rd meeting of the Ministers' Deputies), available at: https://search.coe.int/cm/pages/result_details.aspx?objectId=09000016809e1154.

collection, curation and processing) and domain-specific deployment of algorithmic systems.

A comprehensive approach tackling these four different areas will ensure a robust human rights-based approach. As a result of this study, it is also suggested that the GEC and the CDADI develop, based on those four avenues and through a dedicated Committee of Experts, a Committee of Ministers Recommendation on the impact of artificial intelligence systems, the discriminatory risks they cause, and exploring their potential for promoting equality, including gender equality. This Recommendation will aim to fill the general framework set by the aforementioned convention with regard to the principle of equality.



First avenue: Prevention, transparency and accountability

*This Study has demonstrated a need to shift the prevailing paradigm in anti-discrimination law. Whereas anti-discrimination obligations currently focus on redressing harms done, this Study shows that tackling algorithmic discrimination requires complementing existing ex-post control and redress mechanisms by introducing or strengthening legal obligations pertaining to preventive measures and ex-ante safeguards. A preventive approach is necessary due to the foreseeability and scale of algorithmic bias and resulting harms: **where no preventive safeguards are put in place, AI reproduces and amplifies existing patterns of inequality.** Such a shift in the regulatory paradigm is also justified by the new power and information asymmetries that arise between providers, users and subjects of algorithmic decision-making systems to the disadvantage of potential victims of algorithmic discrimination, as well as by the new vulnerabilities that stem from pervasive social sorting powered by predictive analytics. Transparency obligations and accountability mechanisms can serve to mitigate such asymmetries and to empower end-users of algorithmic systems, including potential victims of algorithmic discrimination and those representing their interests*

- 1) Member states are encouraged to **expand the use of positive action measures to tackle algorithmic discrimination** and to **use the concept of positive obligations anchored in ECHR case law to create an obligation for providers and users to reasonably prevent algorithmic discrimination.** Such positive action measures could be modelled on existing positive obligations under the ECHR, EU law and CRPD provisions on “reasonable accommodation” in relation to discrimination on grounds of disability. Neutrality towards algorithmic bias will not prevent algorithmic discrimination. The principle of reasonable accommodation offers a useful legal yardstick to think about preventive measures. It allows scaling the costs of preventive measures to the size and economic power of the user involved. It also puts positive action and the substantive approach to promoting equality, including gender equality, at the centre of legal responses to algorithmic discrimination. Interpretive guidelines concerning Article 14 ECHR endorsing a positive obligation to prevent algorithmic discrimination should also be made available.
- 2) Member states are encouraged to **introduce mandatory discrimination risk and equality impact assessments** throughout the lifecycle of algorithmic systems according to their specific uses. *Ex-ante* accountability and justification obligations could, for example, require providers and users of AI systems to perform preliminary discrimination risk and equality impact assessments independently or as part of a broader human rights impact

assessment. The results of discrimination risk and equality impact assessments could determine whether a given algorithmic system can enter the market and be used. Such a “pre-approval system” obliges providers and users of potentially risky systems to take preventive measures and possibly submit their systems to certification or a licensing.²¹⁷ Further scrutiny should be applied to systems deemed as posing serious risks of discrimination as explained below. The legal principle of proportionality could help scale these obligations to the size, capacity and economic power of providers and users. Such risk and impact assessments should be made public and easily accessible. They should assess the potential discriminatory impact of algorithmic systems throughout their entire life cycle and across the range of grounds protected under the European Convention on Human Rights. In the framework of HUDERIA, the Human Rights, Democracy and Rule of Law Impact Assessment proposed by the Council of Europe, the main elements of such assessment could be:

1. *Risk Identification*: Identification of relevant risks for equality, including gender equality, and non-discrimination;
2. *Impact Assessment*: Assessment of the impact, taking into account the likelihood and severity of the effects on equality rights;
3. *Governance Assessment*: Assessment of the roles and responsibilities of duty-bearers, rights holders and stakeholders in implementing and governing the mechanisms to mitigate the impact;
4. *Mitigation and Evaluation*: Identification of suitable mitigation measures and ensuring a continuous evaluation.

Member states are encouraged to render such assessments legally binding or to create a strong incentive by making such impact assessment an element to be considered by judges when called upon to assess claims of algorithmic discrimination and whether positive obligations pertaining to equal treatment have been met. So far, some jurisdictions have already proposed legislation that would implement algorithmic impact assessments as a tool to bring accountability to the algorithmic systems increasingly used in everyday life.²¹⁸ Despite this heightened focus on impact assessments as an

217. See G. Malgieri and F. Pasquale, ‘From Transparency to Justification: Toward Ex Ante Accountability for AI’ (2022) Brooklyn Law School, Legal Studies Paper N. 712.

218. The Netherlands has recently made such human rights impact assessments mandatory for public institutions. See the motion here: <https://www.tweedekamer.nl/kamerstukken/moties/detail?id=2022Z06024&did=2022D12329>.

algorithmic governance mechanism, no truly accountable process for conducting such assessments has been standardised.²¹⁹

3) Certification mechanisms could be used in addition to ensure that biases have been mitigated and risks of discrimination eliminated as far as possible for well-defined uses. Such *ex-ante* prevention measures could be used to assert the non-discriminatory nature of domain-specific systems. Certification should enjoy a degree of publicity.²²⁰ For example, it should be made accessible and intelligible to users and end-users of algorithmic systems (in view of possible legal defence and risk assessments for certain applications) as well as to public institutions (e.g., equality bodies, national human rights institutions, data protection officers...) and to CSOs with a legitimate interest (as defined in the EU equality directives and national equality law). The Recommendation CM/Rec(2020)1 on the human rights impacts of algorithmic systems published by the Council of Europe in 2020 could be used as inspiration here. It provides that '[c]ertification schemes based on regional and international standards should be designed and applied to guarantee the provenance and quality of datasets and models' and that they 'should also form part of procurement processes and should be informed by, and compliant with, regulatory frameworks that ban certain uses of algorithmic systems'.²²¹

4) Member states are encouraged to investigate the relationship between accountability, transparency and trade secrets law as it pertains to AI. For an algorithm to be explainable it needs to have a degree of accessibility, whether by internal or external auditors, a regulator, or a tribunal. However, a company's own algorithm may also be covered by trade secrets legislation. There are interesting developments in this sense thanks to the emergence of Secure Multi Party Computation that may enable an

219. Yet, methodologies have been developed in scholarly research, see also Mantelero, Alessandro and Esposito, Samantha, An Evidence-Based Methodology for Human Rights Impact Assessment (HRIA) in the Development of AI Data-Intensive Systems (March 22, 2021). Computer Law & Security Review, 2021 and Mantelero, Alessandro. "Human Rights Impact Assessment and AI." Beyond Data. TMC Asser Press, The Hague, 2022. 45-91. Some examples such as the Dutch "Impact Assessment Fundamental rights and algorithms" can serve as a source of inspiration, see Janneke Gerards, Mirko Tobias Schäfer, Arthur Vankan, Iris Muis, "Impact Assessment Fundamental rights and algorithms", Ministry of the Interior and Kingdom Relations (2022) available at: <https://www.government.nl/binaries/government/documenten/reports/2022/03/31/impact-assessment-fundamental-rights-and-algorithms/Fundamental+Rights+and+Algorithms+Impact+Assessment.pdf> (last accessed 22 July 2022).

220. Publicity should also help mitigate the flows of certification, e.g., in relation to conflicts of interests.

221. Recommendation CM/Rec(2020)1 of the Committee of Ministers to member States on the human rights impacts of algorithmic systems.

interrogation of AI without access to the actual code. Further research should be encouraged in this area, so that the necessary scrutiny in the context of discrimination claims may be calibrated to account for the need to preserve industrial secrecy.

5) **Member states are encouraged to consider establishing legal obligations for users of AI systems to publish statistical data** that can allow interested parties to assess the discriminatory effects of a given system. The data required should include disaggregation on how the system affects the groups protected by the ECHR anti-discrimination provisions, including potential intersectional discrimination. Data should also be published regarding the functioning of algorithmic systems including, when possible, the decision criteria used, and, where relevant, information about the training and validation data used and its processing. Transparency obligations should be balanced against trade secrets and data protection rules. In particular, member states are encouraged to establish legal requirements for users of algorithmic systems to provide consumers and public oversight bodies with a meaningful explanation upon request, especially as part of ongoing legal proceedings. This could be achieved, for instance, by ensuring an effective a right to accessible and intelligible information.

6) Member states are encouraged to **introduce mechanisms for transparency**, which may include annual reporting of AI use in local authorities and government; an obligation for companies to report on responsible AI use through their annual reporting processes, including through the Environmental, Social and Governance (ESG) requirements; and the use of registers of algorithms. The city of Amsterdam, for example, introduced such a register to provide an overview of the artificial intelligence systems and algorithms used by the city.²²²

7) Prevention, transparency and accountability measures should be consolidated by member states in a **comprehensive Action Plan on AI and Equality** that can inform the broader public of ongoing initiatives and guide concrete stakeholder efforts. The Action Plan of member states should help devise a comprehensive 'equality by design' strategy to mainstream equality, anti-discrimination and gender perspectives in the development of AI and algorithmic systems. In addition, member states are encouraged to adopt policies to facilitate the collection of equality data to support the assessment of algorithmic systems in relation to their discriminatory impact. Equality data should also take into account intersectional discrimination.

222. Amsterdam Algorithmic Register, available at: <https://algorithmerregister.amsterdam.nl/en/ai-register/>

Second avenue: Access to justice and legal redress mechanisms

As this Study demonstrates, the deployment of AI and algorithmic systems upsets some of the traditional power equilibria on which the central elements of non-discrimination law rest. In particular, new information asymmetries undermine some of the most fundamental assumptions regarding victims' access to evidence, the ability of victims to establish, and of defendants to rebut, prima facie discrimination, the allocation of liability, the identification of causal relationships in the emergence of discrimination, and the assessment of proportionality and justifications by judges. This endangers access to justice and the integrity and effectiveness of existing legal redress for algorithmic discrimination. Because it is unclear how citizens and consumers will be able to bring a claim forward and what avenues they can pursue for such purpose, the following steps are encouraged:

1) Member states are encouraged to **facilitate access to justice by establishing public supervision mechanisms and developing collective action routes for investigation and redress of algorithmic discrimination**. Public bodies and organisations including equality bodies, national human rights institutions, ombudspersons and data protection agencies should be explicitly mandated to monitor algorithmic discrimination in member states (e.g., through audits of algorithmic systems); to disseminate information and raise awareness among the public; to investigate potential cases of algorithmic discrimination (e.g., to assess compliance with existing laws through testing procedures); to support victims, including with free legal counselling and aid; and to redress algorithmic discrimination either through the power to issue binding opinions and/or to represent victims and/or to intervene in court proceedings. This requires staff training, capacity-building and adequate funding. Existing structures (e.g., equality bodies under EU equality law) should be reinforced for this purpose and co-operation and synergies should be exploited to use existing know-how and competences. As part of this, it will be important to take into account the two proposals for directives on standards for equality bodies that are currently being negotiated in the EU.²²³ If adopted, these directives will strengthen the role that equality bodies can play in monitoring and investigating cases of algorithmic discrimination and

223. Proposal for a Directive of the European Parliament and of the Council on standards for equality bodies in the field of equal treatment and equal opportunities between women and men in matters of employment and occupation COM(2022) 688 final and Proposal for a Council Directive on standards for equality bodies in the field of equal treatment between persons in matters of social security and in the access to and supply of goods and services COM(2022) 689 final.

in supporting victims. For instance, the proposed directives would entrust further investigative powers to equality bodies, that will also be able to issue opinions or binding decisions, act in court cases and offer an alternative dispute resolution mechanism for discrimination claims.

2) **Member states are encouraged to adjust, complement and reinforce the effectiveness of evidence rules to create a fairer, more balanced burden of proof.** Several suggestions could be used as inspiration:

- ▶ Rebuttable presumptions regarding the (lack of) respect for equality and non-discrimination obligations could be explored. For instance, a provisional reversal of the presumption of lawfulness of algorithmic systems could better reflect both the pervasiveness and foreseeability of algorithmic bias and the power shifts described above. This proposal, also aligned with the work of the EU in the framework of the AI Act, could align with the EU's proposal to create a rebuttable presumption of causality towards defendants in case of AI-related damages.²²⁴ Flanked by disclosure obligations, such a presumption would result in alleviating victims' burden of proof when it comes to showing discrimination *prima facie*. It resonates with the proposal for a "distrust by design" framework.²²⁵ Arguably, the foreseeable algorithmic harms and existing power imbalances justify the existence of a rebuttable presumption of algorithmic bias.
- ▶ Providers and users of algorithmic systems would be able to prevent such presumptions of algorithmic bias from arising where they can show that they have taken meaningful and sufficient preventive measures. Such measures would include, for example, discrimination risk and equality impact assessments, and certification (see above).
- ▶ Where such preventive measures have not been taken, a provisional presumption of algorithmic bias could be posited, with the effect of shifting the burden of proof towards the defendant, relying on existing rules on reversing the burden of proof. Users and/or providers of algorithmic systems would then need to show that the system complies with anti-discrimination law requirements.
- ▶ Such *ex-ante* accountability mechanisms would also assist judges and equality bodies with a decision-making function in applying the

224. As explained in Section 2.III. the recent proposal for an AI Liability Directive could offer some inspiration for easing the burden of proof in cases of algorithmic discrimination, including rebuttable presumptions and disclosure obligations.

225. Malgieri G and Pasquale F, 'From Transparency to Justification: Toward Ex Ante Accountability for AI' (2022) Brooklyn Law School, Legal Studies Paper N. 712.

proportionality test in cases of discrimination and to evaluate the adequacy of preventive safeguards taken. This mechanism would create incentives for users and providers of AI to translate technical choices into legally intelligible information which can then be used to assess whether a given system fulfils a legitimate aim, and whether the means employed are reasonably proportionate to the aim pursued.

- ▶ When it comes to liability, member states could explore the concept of negligence to design legal responses where no preventive measures against algorithmic bias have been taken by providers and/or users of a discriminatory system. In the absence of legal personhood of AI systems, users and providers should strategically be held liable for algorithmic discrimination, even where it has been ‘autonomously’ produced by the discriminatory system. Failure to reasonably prevent algorithmic discrimination could give rise to a rebuttable presumption as set out earlier.
- ▶ Liability for algorithmic discrimination should yield clear requirements to remove any discriminatory impact in given uses and to compensate for (moral and material) damage suffered because of algorithmic discrimination. Clear time limits should be adopted for removing any discriminatory impact from AI systems. Publication obligations could also support the enforcement of liability rules, for instance with mandatory information to be provided to the public by any private undertaking found liable for algorithmic discrimination.

3) **Encourage co-operation between regulatory bodies and agencies.**

For example, financial services regulators should be enabled to share work with data protection authorities, especially as regards the use of personal data in the context of allocating loans or investigating potential bias. Furthermore, regulators should research the use of privacy-enhancing technologies and standardised audits mechanisms. In the same vein, member states are encouraged to facilitate co-operation between data protection authorities, equality bodies, national human rights institutions and consumer protection agencies.

4) In line with the CAI draft framework, which pays attention not only to discrimination but also to vulnerability, members states are encouraged to **investigate the new forms of ‘algorithmic’ vulnerability that emerge with the use of AI systems** and to **encourage research into methodologies for legal protection against** such vulnerability. The pervasive deployment of AI systems generates extreme forms of social sorting and differentiation that question the boundaries of discrimination law. Such social sorting creates new forms of social vulnerability, for example by subjecting certain

algorithmic groups to systematically worse economic conditions or exclusion from certain goods and services. These new forms of social differentiation are likely to create “emergent discrimination” as well as deeply entrenched socio-economic inequalities.²²⁶ In addition, it is important to recognise the new conditions of “inferiority, dependency, and subjugation” experienced by vulnerable individuals in the context of data processing.

5) Member states are encouraged to clarify that the prohibition of discrimination in Art.14 covers **intersectional discrimination and discrimination by proxy**, two forms of discriminatory harms that algorithmic systems are most likely to generate. At Council of Europe level, this could be supported by developing recommendations and guidelines on the interpretation of Art. 14 ECHR.

6) In line with the work of the CAI, member states are encouraged to explore how **consumer protection law could be used to complement anti-discrimination**, for instance by facilitating access to information, prohibiting certain features in algorithmic systems under the notion of abusive clauses, etc.

Third avenue: Diversity, inclusion, representation and participation

This study underlines the need to avoid techno-centrism and solutionism and to focus on discrimination, its social components and its roots rather than on technical bias.²²⁷ A core component of mainstreaming equality in AI is to ensure diversity and inclusion through the representation and participation of women and discriminated groups in relevant professional communities and through dedicated training. In particular, educational institutions and businesses should be required to foster a culture of openness, inclusion and diversity through positive action measures aimed to widen access to and success in the AI professional sectors. Therefore, member states are encouraged to take the following steps:

1) Member states should **identify and support positive action measures**, including measures on increasing women’s participation and diversifying

226. See Matthias Leese, The new profiling: Algorithms, black boxes, and the failure of anti-discriminatory safeguards in the European Union, 45 SECURITY DIALOGUE 494–511, 501 (2014); Monique Mann & Tobias Matzner, Challenging algorithmic profiling: The limits of data protection and anti-discrimination in responding to emergent discrimination, 6 BIG DATA & SOCIETY, 5–6 (2019).

227. For instance, Mantelero and Fanucci call for “adopting a broader perspective on the AI industry by considering its entire supply chain from a human rights standpoint”. See Mantelero, Alessandro and Fanucci, Francesca, Great Ambitions. The International Debate on AI Regulation and Human Rights in the Prism of the Council of Europe’s CAHAI (April 4, 2022). Philip Czech et al. (eds). European Yearbook on Human Rights 2022 (Intersentia, Forthcoming).

professional communities, and **should actively enforce positive obligations to promote equality in AI**. Member states should consider making **some forms of positive action legally binding** in the context of tackling algorithmic discrimination. Member states are particularly encouraged to **mainstream positive action in professional sectors involved in the non-discriminatory development, risk assessment and deployment of AI systems**. Positive action measures should apply to the AI professional community, for example through supporting the long-term integration of women and girls (for instance with measures supporting work-life balance), persons with diverse ethnic, religious and linguistic backgrounds, LGBTI persons and other protected and disadvantaged groups in relevant educational programmes and job sectors.²²⁸ Several policy options ranging from awareness-raising to financial support, dedicated training and quotas can be envisaged for this purpose. Such measures could include, for instance:

- ▶ Hiring and promotion policies aiming to achieve a balanced representation of people exposed to discrimination in highest-level positions and support programmes to help position them for success, including a requirement to hire and promote women for at least 50% of the highest-level positions in tech sectors;
- ▶ Quotas and scholarship programmes to support the representation of women and girls and persons exposed to discrimination in STEM studies;
- ▶ Sectoral and/or company-wide policies on addressing stereotypes, discrimination, harassment and violence against women and protected groups in the workplace, including workplace policies fostering inclusion and re-integration following career breaks;
- ▶ Training programmes to raise the awareness of the industry workforce of the discriminatory effects of AI and strategies to prevent them.

2) **Positive obligations to promote equality should provide the legal basis to ensure that AI and algorithmic systems are developed with equality promotion at their core**. Enforcing such obligations more actively could yield a paradigm shift in relevant professional communities and could ensure that AI and algorithmic products are designed and developed with

228. For examples of supporting studies, see e.g., EIGE, Study and work in the EU: Set apart by gender (2018) available at: <https://eige.europa.eu/publications/study-and-work-eu-set-apart-gender-report>; EIGE, Gender Equality Index 2020 Digitalisation and the future of work (2020) available at: <https://eige.europa.eu/publications/gender-equality-index-2020-digitalisation-and-future-work> and EIGE, Artificial intelligence, platform work and gender equality (2020) available at: <https://eige.europa.eu/publications/artificial-intelligence-platform-work-and-gender-equality>.

equality in mind. For instance, member states could establish **procurement rules** that include admissibility requirements for diversity in the professional teams responsible for the development of AI systems as well as to demonstrate how the rule of law has been incorporated into the AI system used by public authorities.

3) Positive obligations to promote equality could also translate into a requirement for companies of the AI sector to **develop and implement an equality strategy** covering the groups protected under Article 14 ECHR and Article 1 Protocol 12 ECHR. Member states could also encourage or require AI sector companies to **appoint an officer in charge of:**

- ▶ supervising the enforcement of this strategy;
- ▶ **forging dialogue between legal and technical teams** to facilitate the introduction of legal requirements in the design of AI systems; and
- ▶ **co-operating with regulatory and enforcement authorities to demonstrate the compliance** of the systems developed by the company with non-discrimination law or, where necessary, make changes to ensure compliance or withdraw these systems from the market.²²⁹

Fourth avenue: Democratic participation, public awareness-raising and capacity-building

There is no doubt that AI is reshaping our societies, and, alongside ex-ante requirements, regulatory enforcement and ex post assessments, it is crucial for member states to invest in educating citizens and consumers to fully empower their digital citizenship. Therefore:

1) Member states are encouraged to introduce a **right to information on algorithmic mediation** in the context of discrimination complaints or claims. A requirement should be placed on all organisations to inform users as to whether they are interacting with a human or a machine. Users should also be informed about how decisions are made and how they can be challenged.

2) Member states should encourage the **rolling out of digital literacy programmes**, especially in a context-aware manner, to improve awareness of rights to equality, including gender equality, and non-discrimination in the context of AI applications.²³⁰ Member states should also encourage a **culture fostering collective bargaining in relation to digitalization of the**

229. In this sense, see Yeung, Karen and Harkens, Adam, How Do 'Technical' Design-Choices Made When Building Algorithmic Decisionmaking Tools for Criminal Justice Authorities Create Constitutional Dangers? (Part 1) (December 7, 2022). Public Law, Forthcoming, p. 3.

230. There is also an issue of access to AI applications, in particular if a paywall restricts their use.

workplace and the participation of stakeholders in algorithmic management decisions. For example, in a labour context, the Italian trade union CGIL (*Confederazione Generale Italiana del Lavoro*) proposed to review and negotiate the ways in which algorithmic systems are involved in the organisation of work and working processes as part of trade union negotiations).²³¹

3) Member states are encouraged to strengthen legal requirements on democratic participation in standard-setting. Standards will play a crucial role in the exposure of AI systems to market forces. However, the standard setting process is taking place far from public scrutiny and involvement of organisations as well as individuals, especially the most vulnerable. Therefore, member states are encouraged to take the following action:

- ▶ **Identify good practices** for the democratization of the AI standard setting process.
- ▶ **Establish legal requirements for democratic participation** in the AI standards development processes. In particular, as AI standardization plays an increasingly prominent role, including for equality and non-discrimination, there is an inherent risk of standards being used to define and interpret legal requirements which impact upon human rights, especially as this is happening through private standardization bodies with limited participation from civic society. For this reason, in order to ensure the representation of women, groups affected by discrimination and of equality experts, member states are encouraged to facilitate the participation of these groups in standard-setting processes, including NGOs with a legitimate interest. Public consultations and public oversight are required, as standard-setting is at risk of veering into the areas of public policy and law, which may require a level of interpretation, such as bias in data.

4) Capacity-building should include investment into interdisciplinary research on non-discriminatory algorithms and into strategies to protect equality in the use of algorithmic systems. Nevertheless, recent research has confirmed that preventing algorithmic discrimination by eliminating bias, for example by debiasing datasets, is a highly unlikely prospect, both because of the complexity of algorithmic bias and because of the evolving nature of machine-learning systems, especially when put to use in dynamic

231. See Daniele Carchidi, 'Contrattare per governare gli impatti della digitalizzazione sul mondo del lavoro: il caso afiniti' (2022) SLC-CGIL available at: <https://www.slc-cgil.it/notizie-tlc-ed-emittenza/3791-afiniti-un-caso-riuscito-di-contrattazione-dell-algoritmo.html>.

social contexts.²³² While the existence of bias is foreseeable, predicting and checking for all possible biases and any ensuing discrimination is an impossible task.²³³ Debiasing can, therefore, only represent one aspect of such preventive measures and should be **complemented by measures addressing the societal roots of discrimination.**

232. See Balayn A and Gürses S, Beyond Debiasing: Regulating AI and its inequalities (European Digital Rights 2021).

233. Ibid.

www.coe.int

The Council of Europe is the continent's leading human rights organisation. It comprises 46 member states, including all members of the European Union. All Council of Europe member states have signed up to the European Convention on Human Rights, a treaty designed to protect human rights, democracy and the rule of law. The European Court of Human Rights oversees the implementation of the Convention in the member states.

COUNCIL OF EUROPE



CONSEIL DE L'EUROPE