

# Escaping the Impossibility of Fairness: From Formal to Substantive Algorithmic Fairness

Ben Green, University of Michigan ([bzgreen@umich.edu](mailto:bzgreen@umich.edu))

## Abstract

As governments embrace algorithms, the burgeoning field of algorithmic fairness provides an influential methodology for promoting equality-enhancing reforms. However, even algorithms that satisfy mathematical fairness standards can exacerbate oppression, causing critics to call for the field to shift its focus from “fairness” to “justice.” Yet any efforts to achieve algorithmic justice in practice are constrained by a fundamental technical limitation: the “impossibility of fairness” (an incompatibility between mathematical definitions of fairness). The impossibility of fairness thus raises a central question about algorithmic fairness: How can computer scientists support equitable policy reforms with algorithms? In this article, I argue that promoting justice with algorithms requires reforming the methodology of algorithmic fairness. First, I diagnose why the current methodology for algorithmic fairness—which I call “formal algorithmic fairness”—is flawed. I demonstrate that the problems of algorithmic fairness—including the impossibility of fairness—result from the methodology of the field, which restricts analysis to isolated decision-making procedures. Second, I draw on theories of substantive equality from law and philosophy to propose an alternative methodology: “substantive algorithmic fairness.” Because substantive algorithmic fairness takes a more expansive scope to fairness, it enables an escape from the impossibility of fairness and provides a rigorous guide for alleviating injustice with algorithms. In sum, substantive algorithmic fairness presents a new direction for algorithmic fairness: away from formal mathematical models of “fairness” and toward substantive evaluations of how algorithms can (and cannot) promote justice.

## 1 INTRODUCTION

### 1.1 Algorithmic Fairness and Its Discontents

Machine learning algorithms have become central components in many efforts to promote equitable public policy. In the face of widespread concerns about discriminatory institutions and decision-making processes, many policymakers and scholars praise algorithms as critical tools for equality-enhancing reforms (Arnold Ventures, 2019; Eubanks, 2018; Harris & Paul, 2017; Porrino, 2017). To policymakers, policy advocates, and scholars across multiple fields, algorithms overcome the cognitive limits and social biases of human decision-makers, enabling more objective and fair decisions (Arnold Ventures, 2019; Harris & Paul, 2017; Kleinberg et al., 2019; Miller, 2018; Sunstein, 2019). Thus, for instance, in light of concerns about the biases of judges, many court systems in the United States have adopted pretrial risk assessments as a central component of criminal justice reforms (Green, 2020; Koepke & Robinson, 2018; Porrino, 2017).

Undergirding these reform efforts is the burgeoning field of algorithmic fairness. Grounded primarily in computer science, algorithmic fairness applies the tools of algorithm design and analysis—in particular, an emphasis on formal mathematical reasoning (Green & Viljoen, 2020)—to fairness. The central components of algorithmic fairness are developing mathematical definitions of fair decision-making (Barocas et al., 2019), optimizing algorithms for these definitions (Feldman et al., 2015; Hardt et al., 2016), and auditing algorithms for violations of these definitions (Angwin et al., 2016; Obermeyer et al., 2019; Raji & Buolamwini, 2019).

In the context of policy reform efforts, algorithmic fairness is often employed to determine whether an algorithm is “fair” and, therefore, appropriate to use for decision-making. For instance, in settings such as pretrial adjudication and child welfare, debates about whether to employ algorithms hinge on evaluations of algorithmic fairness (Angwin et al., 2016; Chouldechova et al., 2018; Dieterich et al., 2016; Eubanks, 2018). Similarly, regulation of government algorithms often calls for evaluations that test algorithms for biases (California Legislature, 2021; European Commission, 2021; Government of Canada, 2021).

Yet as algorithmic fairness has risen in prominence, critical scholars have highlighted several concerns. Efforts to formulate mathematical definitions of fairness overlook the contextual and philosophical meanings of fairness (Binns, 2018; Green & Hu, 2018; Jacobs & Wallach, 2021; Selbst et al., 2019). Algorithmic fairness focuses on bad actors, individual axes of disadvantage, and a limited set of goods, thus “mirroring some of antidiscrimination discourse’s most problematic tendencies” as a mechanism for achieving equality (Hoffmann, 2019). As a result, there is often a significant gap between mathematical evaluations of fairness and an algorithm’s real-world impacts (Green & Viljoen, 2020). Algorithms that satisfy fairness standards often exacerbate oppression and legitimize unjust institutions (Davis et al., 2021; Green, 2020; Kalluri, 2020; Ochigame, 2020; Ochigame et al., 2018; Powles & Nissenbaum, 2018). In turn, some scholars have called for rejecting the frame of “fairness” altogether, proposing alternative frames of “justice” (Bui & Noble, 2020; Costanza-Chock, 2020; Green, 2018), “equity” (D’Ignazio & Klein, 2020), and “reparation” (Davis et al., 2021).

However, efforts to achieve algorithmic justice in practice are constrained by a fundamental technical limitation: the “impossibility of fairness.” This result reveals that it is impossible for an algorithm to satisfy all desirable mathematical definitions of fair decision-making (Chouldechova, 2017; Kleinberg et al., 2016). An algorithm that is fair along one standard will inevitably be unfair along another standard.<sup>1</sup> Although no mathematical definitions of algorithmic fairness fully encapsulate the philosophical notion of fairness or justice (Binns, 2018; Green & Hu, 2018; Jacobs & Wallach, 2021; Selbst et al., 2019), each captures a normatively desirable principle.

---

<sup>1</sup> I will provide more detail on these fairness definitions and the impossibility of fairness in Section 2.

The impossibility of fairness presents an intractable constraint on efforts to promote equitable public policy using algorithms: any effort to improve decision-making using algorithms will violate at least one normatively desirable fairness principle. This result suggests that the best algorithm developers can do to promote justice in practice is tune algorithms to align with some (limited) fairness definitions at the expense of others (Davis et al., 2021; Kleinberg et al., 2019; Wong, 2020). As one article about algorithmic fairness concludes, “the tradeoff between [...] different kinds of fairness has real bite” and means that “total fairness cannot be achieved” (Berk et al., 2018). Similarly, a proponent of algorithmic justice acknowledges that, because of the legal implications of these tradeoffs, “it is highly unlikely that an algorithmic justice approach will advance” (Costanza-Chock, 2020).

The impossibility of fairness thus raises a central question about algorithmic fairness: How can computer scientists support equitable policy reforms with algorithms in practice, given the impossibility of fairness? In this article, I argue that achieving this goal requires reforming the methodology of algorithmic fairness. This argument involves two tasks. Task 1 is to diagnose why the current methodology for algorithmic fairness is flawed. I demonstrate that the problems of algorithmic fairness—including the impossibility of fairness—result from the dominant methodology of the field, which restricts analysis to isolated decision-making procedures. Task 2 is to develop an alternative approach that operationalizes a social justice orientation into algorithmic fairness. Drawing on theories of substantive equality from law and philosophy, I propose a new methodology for algorithmic fairness that enables an escape from the impossibility of fairness and that, in turn, is better equipped to alleviate injustice. This proposed method provides concrete steps to help reform-minded computer scientists rigorously pursue substantive equality with algorithms.

## **1.2 Article Overview: Methodological Reform**

A methodology is “a body of methods, rules, and postulates employed by a discipline” (Merriam-Webster, 2021). A methodology provides a systematic language for comprehending and reasoning about the world, shaping how practitioners formulate problems and develop solutions to those problems. Problem formulation has both practical and normative stakes (Passi & Barocas, 2019). As philosopher John Dewey writes, “The way in which [a] problem is conceived decides what specific suggestions are entertained and which are dismissed” (Dewey, 1938). An inadequately conceived problem “cause[s] subsequent inquiry to be irrelevant or to go astray;” the remedy is to reformulate the problem (Dewey, 1938). Furthermore, as philosopher Elizabeth Anderson describes, “Sound political theories must be capable of representing normatively relevant political facts. If they can’t represent certain injustices, then they can’t help us identify them. If they can’t represent the causes of certain injustices, then they can’t help us identify solutions” (Anderson, 2009). In sum, if a methodology fails to account for normatively relevant facts and principles, it will generate problem formulations that yield unhelpful or unjust proposals for reform.

Critiques of algorithmic fairness and the impossibility of fairness suggest that algorithmic fairness has been conceived in a deficient manner, leading to myopic suggestions for how to promote fairness in practice. Thus, in the spirit of Dewey and Anderson, this article proposes methodological reforms so that algorithmic fairness can provide a more rigorous guide for promoting justice with algorithms. The article proceeds in Section 2 with background on the impossibility of fairness. Section 3 describes how egalitarian theories of substantive equality can inform work on algorithmic fairness. I then turn to the two primary tasks of this article.

Section 4 takes on task 1: diagnose why the current methodology for algorithmic fairness is flawed. I argue that the flaws of algorithmic fairness result from a significant methodological limitation: algorithmic fairness relies on a narrow frame of analysis restricted to specific decision points, in isolation from the context of those decisions.<sup>2</sup> I call this method “formal algorithmic fairness,” as it aligns with formal equality (which emphasizes equal treatment for individuals based on their attributes or behavior at a particular decision point). Formal algorithmic fairness represents a systematic approach to problem formulation in which fairness is operationalized in terms of isolated decision-making processes. Because formal algorithmic fairness is conceived so narrowly, it provides an ineffective guide for achieving algorithmic justice in practice. Formal algorithmic fairness yields a misguided and techno-centric reform strategy: enhance fairness by optimizing decision-making procedures with algorithms. These algorithmic interventions often exacerbate oppression and are constrained by the impossibility of fairness. Thus, formal algorithmic fairness leaves reform efforts in a bind: it appears that our only options are to pursue superficially “fair” algorithms or to reject algorithmic reforms, leaving the status quo in place.

In light of these flaws, Section 5 takes on task 2: propose an alternative approach to algorithmic fairness that enables a justice-oriented agenda for developing and applying algorithms. I call this method “substantive algorithmic fairness,” as it draws on theories of substantive equality from law and philosophy (which emphasize the identification and reduction of social hierarchies<sup>3</sup>). My goal is not to incorporate substantive equality into a formal mathematical model: this strategy that would fail to provide the necessary methodological shift (Green & Viljoen, 2020). Substantive algorithmic fairness is not a method for creating “substantively fair algorithms.” Instead, following an “algorithmic realist” approach (Green & Viljoen, 2020), my goal is to develop problem formulations that are grounded in legal and philosophical reasoning about substantive equality. In other words, rather than treat fairness as a mathematical attribute of algorithms in isolation,

---

<sup>2</sup> By decision points, I refer to the specific moments in which decisions are made about individuals. Examples include decisions about whether to release or detain pretrial defendants and decisions about whether to admit or reject college applicants.

<sup>3</sup> Social hierarchies refer to “social relation[s] of rank ordering, typically on a group or categorical basis,” that generate disparities in social and material resources (MacKinnon, 2011).

substantive algorithmic fairness considers the impacts of algorithms in relation to the social and political context.

Because substantive algorithmic fairness expands the frame of analysis beyond isolated decision points, it suggests reform strategies that escape from the impossibility of fairness and, in turn, can promote algorithmic justice in practice. Substantive algorithmic fairness presents a three-step strategy for promoting substantive equality with algorithms: 1) diagnose the substance of the inequalities in question, 2) identify what reforms can remediate the substantive inequalities, and 3) consider whether algorithms can enhance the desired reforms. Substantive algorithmic fairness thus suggests reforms beyond either implementing a superficially “fair” algorithm or leaving the status quo in place. In sum, substantive algorithmic fairness presents concrete steps toward a new method for algorithmic fairness: away from formal mathematical models of “fairness” as an attribute of algorithms and toward substantive evaluations of how algorithms can (and cannot) promote justice.

## **2 THE IMPOSSIBILITY OF FAIRNESS**

In May 2016, journalists at ProPublica reported that a risk assessment algorithm used to judge pretrial defendants in Broward County, Florida was “biased against blacks” (Angwin et al., 2016). This algorithm, known as COMPAS, was created by the company Northpointe and is used by many court systems across the United States.<sup>4</sup> Like other pretrial risk assessments, COMPAS predicts the likelihood that pretrial defendants will recidivate; these predictions are presented to judges to inform their decisions to release or detain each defendant until their trial (Green, 2020; Koepke & Robinson, 2018). ProPublica found that, among defendants who were not arrested in the two years after being evaluated, Black defendants were 1.9 times more likely than white defendants to be misclassified by COMPAS as “high risk” (i.e., subjected to false positive predictions) (Angwin et al., 2016).

This report sparked significant debate about the use of COMPAS in pretrial adjudication. Tech critics responded to ProPublica’s article with outrage about racist algorithms (Doctorow, 2016; O’Neil, 2016). However, Northpointe and numerous academics defended COMPAS, arguing that ProPublica had focused on the wrong measure of algorithmic fairness (Corbett-Davies et al., 2017; Dieterich et al., 2016; Flores et al., 2016; Gong, 2016). These groups asserted that the proper standard of fairness is not whether false positive (and false negative) rates are the same for each race. Instead, they argued that the proper standard of fairness is whether risk scores imply the same probability of recidivism for each race. COMPAS satisfied this notion of fairness, suggesting that the tool was fair.

---

<sup>4</sup> COMPAS stands for Correctional Offender Management Profiling for Alternative Sanctions. Northpointe has since been renamed Equivant.

This debate about whether COMPAS is fair concerns two distinct definitions of algorithmic fairness. The first is “separation,” which is satisfied if all groups subject to an algorithm’s predictions experience the same false positive rate and the same false negative rate.<sup>5</sup> Separation expresses the idea that people who exhibit the same outcome should be treated similarly. ProPublica argued that COMPAS is biased because it violates separation: Black non-recidivists are more likely to be labeled “high risk” than white non-recidivists (Angwin et al., 2016).

The second notion of algorithmic fairness is “sufficiency,” which is satisfied if, among those who receive a particular prediction, all groups exhibit the outcome being predicted at the same rate.<sup>6</sup> Sufficiency expresses the idea that people who are equally likely to exhibit the behavior of interest should be treated similarly. Northpointe and others argued that COMPAS is fair because it satisfies sufficiency: the label of “high risk” signifies a similar probability of recidivism for both Black and white defendants (Corbett-Davies et al., 2017; Dieterich et al., 2016; Flores et al., 2016; Gong, 2016). Sufficiency is the most widely used notion of algorithmic fairness, particularly because machine learning models typically satisfy this principle by default (Barocas et al., 2019).

The COMPAS debate raised a fundamental question for algorithmic fairness: can an algorithm simultaneously satisfy both separation and sufficiency? As computer scientists soon discovered, the answer is no: there is an inevitable tension between these definitions of fairness (Angwin & Larson, 2016; Barocas et al., 2019; Chouldechova, 2017; Kleinberg et al., 2016). This result is known as the “impossibility of fairness.” The only exceptions to the impossibility of fairness involve two exceedingly rare scenarios: the algorithm makes predictions with perfect accuracy, or all groups exhibit the outcome being predicted at the same “base rate” (Kleinberg et al., 2016). Thus, for instance, a pretrial risk assessment will necessarily either misclassify Black and white defendants as recidivists at different rates (violating separation) or yield different predictions for Black and white defendants who are equally likely to recidivate (violating sufficiency).

The impossibility of fairness reflects a harsh and intractable dilemma facing efforts to promote equality using algorithms (Berk et al., 2018). This dilemma is especially troubling in public policy, where algorithms are typically adopted to enhance the fairness of discrete decision-making processes. In these settings, the statistical fairness measures in tension are salient and often grounded by law. The impossibility of fairness raises a particular challenge for proponents of algorithmic justice: because their proposals involve violating sufficiency in favor of alternate measures (Costanza-Chock, 2020; Davis et al., 2021), such attempts would generally be barred by antidiscrimination law (Corbett-Davies et al., 2017; Costanza-Chock, 2020; Hellman, 2020).

---

<sup>5</sup> Separation is aligned with fairness criteria such as error rate balance and balance for the positive/negative class.

<sup>6</sup> Sufficiency is aligned with fairness criteria such as calibration and predictive parity.

Work on algorithmic fairness operates within the constraints posed by the impossibility of fairness, accepting that the best we can do is to choose a single fairness definition (at the expense of others) or to rigorously balance the tradeoff between multiple definitions (Berk et al., 2018; Costanza-Chock, 2020; Davis et al., 2021; Kleinberg et al., 2019; Wong, 2020). Yet as I will describe (using pretrial risk assessments as a case study), both of these responses lead to narrow reforms that uphold unjust social conditions and institutions. Developing a positive agenda for algorithmic justice requires finding a way to develop and apply algorithms without confronting the impossibility of fairness.

### **3 LESSONS FROM EGALITARIAN THEORY**

Identifying a path for algorithmic justice requires diagnosing why the current methodology for algorithmic fairness is limited (task 1) and developing an alternative methodology that better promotes justice (task 2). In order to inform this evolution, I turn to egalitarian theory. Broadly speaking, “Egalitarian doctrines tend to rest on a background idea that all human persons are equal in fundamental worth or moral status” (Arneson, 2013). Although fairness and equality are complex and contested concepts, both share a central concern with comparing the treatment or conditions across individuals or groups, emphasizing the normative value of some form of parity (Arneson, 2013; Gosepath, 2021; Minow, 2021). Indeed, many definitions of algorithmic fairness make explicit reference to equality (Barocas et al., 2019; Berk et al., 2018). Furthermore, egalitarian scholars have confronted many questions that overlap with central debates in algorithmic fairness (Binns, 2018).

#### **3.1 Formal and Substantive Equality**

Reforming algorithmic fairness requires first understanding why the current method of algorithmic fairness leads to injustice. Egalitarian debates between “formal” and “substantive” equality shed light on this methodological deficit and suggest an alternative approach. Just as algorithmic fairness confronts the limits of narrow formulations of fairness, egalitarian theorists have confronted similar limits of narrow formulations of equality. In response, some egalitarian thinkers have devised more expansive formulations that provide a better guide for ameliorating oppression.

A central tension in egalitarian theory is between “formal” and “substantive” equality. Formal equality asserts, “When two persons have equal status in at least one normatively relevant respect, they must be treated equally with regard in this respect. This is the generally accepted *formal* equality principle that Aristotle articulated [...]: ‘treat like cases as like’” (Gosepath, 2021). In practice, formal equality typically refers to a “fair contest” in which everyone is judged according to the same standard, based only on their characteristics at the moment of decision-making (Fishkin, 2014). In the United States, disparate treatment law is grounded in notions of formal equality, attempting to ensure that people are not treated differently based on protected attributes such as race and gender.

Despite being widely adopted, formal equality suffers from a methodological limitation. Because formal equality restricts analysis to specific decision points, it cannot account for the inequalities that often surround those decision points. Formal equality is therefore prone to reproducing existing patterns of injustice. For instance, a formal equality approach to college admissions would evaluate all applicants based solely on their academic qualifications (e.g., grades and test scores). As long as applicants with similar qualifications are treated similarly, formal equality would be satisfied. Yet because of racial inequalities in educational opportunities (EdBuild, 2019), evaluating all students according to a uniform standard would perpetuate racial hierarchy. Thus, although a formal approach may be sufficient in an equitable society, it “would make no sense at all in a society in which identifiable groups had actually been treated differently historically and in which the effects of this difference in treatment continued into the present” (Crenshaw, 1988).

The limits of formal equality have led many scholars to develop an alternative approach: substantive equality. This approach “repudiate[s] the Aristotelian ‘likes alike, unlikes unlike’ approach [...] and replaces it with a substantive test of historical disadvantage” (MacKinnon, 2011). “Its core insight is that inequality, substantively speaking, is always a social relation of rank ordering, typically on a group or categorical basis,” that leads to both material and dignitary inequalities (MacKinnon, 2011). In other words, “hierarchy identifies the substance of substantive equality” (MacKinnon, 2016). Following this reasoning, substantive equality envisions a world free from social hierarchy (MacKinnon, 2011, 2016). In the United States, disparate impact law is grounded in notions of substantive equality (albeit partially (MacKinnon, 2011, 2016)), attempting to ensure that formally neutral rules do not disproportionately burden historically marginalized groups.

Substantive equality provides the methodological capacity to identify and ameliorate social hierarchies. In contrast to formal equality, substantive equality relies on a broad frame of analysis that evaluates decisions in light of social hierarchies. When confronted with instances of inequality, “A substantive equality approach [...] begins by asking, what is the substance of this particular inequality, and are these facts an instance of that substance?”, emphasizing that “it is the hierarchy itself that defines the core inequality problem” (MacKinnon, 2011). For instance, substantive equality recognizes that racial disparities in college admissions reflect a pervasive racial hierarchy in educational and other opportunities. It therefore rejects the formal equality approach to college admissions. Rather than aiming to evaluate all students according to a uniform standard, substantive equality calls for policies that acknowledge this racial hierarchy (such as affirmative action) and that aim to redress this hierarchy (such as improving educational resources in minority school districts).

As Section 4 will describe, the current approach to algorithmic fairness—which I call “formal algorithmic fairness”—is grounded in formal equality and shares many of formal equality’s limits. This analysis suggests the need for an alternative approach grounded in substantive equality—“substantive algorithmic fairness”—which I present in Section 5.

### **3.2 Substantive Approaches to Escaping Equality Dilemmas**

The second task in reforming algorithmic fairness is to develop an alternative methodology for achieving algorithmic justice in practice. Egalitarian scholarship about how to respond to equality dilemmas suggests strategies for promoting substantive equality without being impeded by the impossibility of fairness. Just as algorithmic fairness confronts the impossibility of fairness, egalitarian theorists have confronted similar tensions between notions of equality. In response, some egalitarian thinkers have devised reform strategies that break free from these dilemmas.

In order to glean insights about how algorithmic fairness can escape the impossibility of fairness, I turn to three complementary substantive equality approaches for analyzing and escaping from equality dilemmas:

- In developing her theory of “democratic equality,” philosopher Elizabeth Anderson responds to a “dilemma” that arises in luck egalitarianism (Anderson, 1999).<sup>7</sup> On the one hand, not providing aid to the disadvantaged means blaming individuals for their misfortune. On the other hand, providing special treatment to individuals on account of their inferiority means expressing contempt for the disadvantaged.
- In developing her “social-relations approach” to equality, legal scholar Martha Minow engages with the “dilemma of difference” that arises in legal efforts to deal with differences between individuals (Minow, 1991). On the one hand, giving similar treatment to everyone regardless of their circumstances can “freeze in place the past consequences of differences.” On the other hand, giving special treatment to those deemed “different” risks entrenching and stigmatizing that difference.
- In developing his theory of “opportunity pluralism,” legal scholar Joseph Fishkin addresses the “zero-sum struggles” that arise in efforts to promote equal opportunity (Fishkin, 2014). On the one hand, judging people for an opportunity based solely on their performance or attributes at a particular moment in time (i.e., a “fair contest”) perpetuates inequalities. On the other hand, even approaches that attempt to account for existing inequalities (such as Rawlsian equal opportunity and luck egalitarianism) fail to create a truly level playing field and prompt “extraordinarily contentious” debates.<sup>8</sup>

---

<sup>7</sup> Luck egalitarianism advocates compensating people for inequalities that result from misfortunate but not inequalities that result from choice (Anderson, 1999; Arneson, 2013).

<sup>8</sup> Each of these scholars are aligned with relational egalitarianism, which asserts that “people should relate to one another as equals or should enjoy the same fundamental status” (Arneson, 2013). Although Fishkin is the least explicitly focused on relationships, his analysis has strong overlaps with relational egalitarianism.

The equality dilemmas presented by Anderson, Minow, and Fishkin resemble the impossibility of fairness. Each scholar presents a dilemma between treating everyone following a uniform standard (akin to sufficiency) and providing special treatment to the disadvantaged (akin to separation). In all of these cases, efforts to promote equality are impaired by a seemingly inescapable, zero-sum tension between notions of equality. If we treat everyone following a uniform standard, we risk reproducing inequality. But if we provide special treatment to the disadvantaged, we might stigmatize the disadvantaged and still fail to achieve greater equality.

In the face of these tradeoffs, it appears difficult—if not impossible—to meaningfully advance equality. As Minow notes, “Dilemmas of difference appear unresolvable” (Minow, 1991). In turn, “decisionmakers may become paralyzed with inaction” (Minow, 1991). At best, decision-makers appear to be left with a zero-sum tradeoff between competing notions of equality. Yet as Fishkin notes, “If [...] zero-sum tradeoffs are the primary tools of equal opportunity policy, then trench warfare is a certainty, and any successes will be incremental” (Fishkin, 2014).

What makes Anderson, Minow, and Fishkin particularly insightful for algorithmic fairness is that they provide methodological accounts of how to escape from these dilemmas. Each scholar reveals that their dilemma is not intractable. Instead, each dilemma only appears intractable if one analyzes inequality through a narrow lens, which restricts the range of possible remedies. Expanding the frame of analysis clarifies the problems of inequality and yields two reform strategies that escape these equality dilemmas.

### *3.2.1 The Relational Response*

The first approach to escaping equality dilemmas is what I call the “relational response”: reform institutions and social norms to reduce social hierarchies. This follows from a substantive analysis highlighting how equality dilemmas are driven by social hierarchies. Noting that the goal of egalitarianism is “to end oppression, which by definition is socially imposed,” Anderson expands the analysis of equality from distributions (of both tangible and intangible goods) to equality of social relations (Anderson, 1999). From this perspective, the problem of inequality is not merely that some people have more of a particular good than others. A deeper problem is that society imposes disadvantages on individuals who lack certain attributes or abilities (Anderson, 1999; Minow, 1991).

Recognizing social categories as relational (rather than intrinsic to individuals) and social arrangements as political and mutable (rather than neutral and static) yields reforms that “escape or transcend the dilemmas of difference” (Minow, 1991). In other words, the primary task of reform should not be providing special treatment to “different” individuals. Instead, reform should aim to reduce the extent to which superficial differences lead to significant disparities in status and abilities (Minow, 1991). Without social hierarchies, real or perceived differences between

individuals would not lead to different levels of rights or capacities, which in turn would prevent the dilemma between treating everyone the same and providing special treatment.

For instance, the injustice faced by someone who is stigmatized because of their physical appearance is not that they are inherently ugly (indeed, the notion of inherent ugliness should be contested). Instead, “the injustice lies [...] in the social fact that people shun others on account of their appearance” (Anderson, 1999). Oppressive social norms turn a superficial difference between people into one marked by severe disparities in status. This feature of social relations creates a dilemma. Treating everyone the same would leave “ugly” individuals in a subordinate position. However, a remedy such as subsidizing plastic surgery for “ugly” individuals would uphold oppressive beauty norms even if it provides aid for some people.

The relational response suggests a strategy that escapes from this dilemma: alter social norms so that no one is shunned or treated as a second-class citizen due to their appearance. If one’s appearance has no relationship to their social status, then appearance ceases to be a normatively relevant category, such that there is no dilemma between treating people similarly or differently based on how they look. Such reforms may be difficult to achieve (at least in the immediate term), thus necessitating more individualized remedies. Nonetheless, this approach “lets us see how injustices may be better remedied by changing social norms and the structure of public goods than by redistributing resources” (Anderson, 1999).

### 3.2.2 *The Structural Response*

The second approach to escaping equality dilemmas is what I call the “structural response”: reduce the scope and stakes of decisions that exacerbate social hierarchies. This follows from a substantive analysis highlighting how the structure of decisions exacerbates social hierarchies and raises the stakes of equality dilemmas. Fishkin broadens the focus from individual competitions to the entire structure of opportunities. From this perspective, the problem of inequality is not merely that groups face vastly different development opportunities, making it impossible to create fair contests between all individuals. A deeper problem is that opportunities are structured around a small number of “zero-sum, high-stakes competitions,” which Fishkin calls “bottlenecks” (Fishkin, 2014). These competitions typically hinge on attributes that are unequally distributed across groups, compounding existing disadvantage (i.e., oppressed groups are less qualified to succeed in competitions for beneficial opportunities, such as jobs).

Fishkin suggests, “Instead of taking the structure of opportunities as essentially given and focusing on questions of how to prepare and select individuals for the slots within that structure in a fair way, [we should] renovate the structure [of opportunities] itself” (Fishkin, 2014). In other words, the primary task of reform should not be helping some disadvantaged individuals receive favorable decisions through special treatment. Instead, reform should aim to limit the extent to which high-

stakes decisions hinge on attributes that are unevenly distributed across social groups due to oppression. Without these bottlenecks, decisions would not as strongly magnify existing inequalities, which in turn would lower the stakes of the dilemma between treating everyone the same and providing special treatment.

For instance, debates about admission to elite US colleges and universities are contentious not only because of inequities in educational resources, but also because admission provides a rare pathway to high social status and material comfort. The significance of college admissions decisions makes disparities in primary and secondary education particularly consequential for determining future life outcomes. These stakes of college admissions creates a dilemma. Evaluating all students according to the same standard would entrench inequalities in primary and secondary education. However, attempts to promote equality through affirmative action are inevitably zero-sum and leave the bottleneck in place.

The structural response provides an escape from this dilemma: lower the stakes of college admissions decisions. Making college admissions less determinative of future life outcomes would reduce the downstream harms of disparities in early educational opportunities. Achieving this goal requires altering the structure of opportunities to create more paths for people to lead comfortable and fulfilling lives without a college degree. By making inequities in primary and secondary education less consequential, these reforms would reduce the dilemma between treating college applicants similarly or differently based on their academic performance.

The relational and structural responses present two concrete substantive equality approaches for dealing with equality dilemmas. As Section 5 will describe, substantive algorithmic fairness applies these substantive equality strategies to the impossibility of fairness. Following the relational and structural responses enables algorithms to escape the impossibility of fairness and to alleviate social hierarchies.

#### **4 FORMAL ALGORITHMIC FAIRNESS: NAVIGATING THE IMPOSSIBILITY OF FAIRNESS**

This section focuses on the first task of reforming algorithmic fairness: diagnosing the current limits of algorithmic fairness as a guide for promoting equitable public policy. I characterize the dominant method of algorithmic fairness as “formal algorithmic fairness.” Akin to formal equality, formal algorithmic fairness limits analysis to the functioning of algorithms at particular decision points. When confronted with concerns about discriminatory decision-making, formal algorithmic fairness formulates the problem in terms only of the inputs and outputs of the decision point in question. As a result, fairness is defined as a technical attribute of algorithms: all of the major definitions of algorithmic fairness are based on the statistical properties of an algorithm’s outputs (Barocas et al., 2019; Berk et al., 2018).

Due to its narrow frame of analysis, formal algorithmic fairness suffers from many of the same methodological limits as formal equality as a guide to achieving equality. To elucidate these limits, I interrogate the two responses to the impossibility of fairness that arise within formal algorithmic fairness. These responses reveal how formal algorithmic fairness yields reforms that appear fair but in practice reproduce injustice. Even the best-case scenario within formal algorithmic fairness provides a meager strategy for promoting equity. All told, the central problem facing algorithmic fairness is not that we lack the appropriate formal definitions of fairness, that data is often biased, or that we cannot achieve sufficient predictive accuracy. The problem is that the method of formal algorithmic fairness restricts analysis to algorithms in isolation, trapping reform efforts within the impossibility of fairness.

#### **4.1 The Fair Contest Response: Reproducing Inequity**

The first formal algorithmic fairness response to the impossibility of fairness is what I call the “fair contest response.” This response defends sufficiency as the proper definition of algorithmic fairness, asserting that fairness entails treating people similarly based solely on each person’s likelihood to exhibit the outcome of interest. On this view, as long as an algorithm satisfies sufficiency, any lack of separation is acceptable—it is the inevitable byproduct of groups exhibiting the outcome in question at different rates. This response applies the logic of a “fair contest,” aiming to evaluate everyone based only on their characteristics at the moment of decision-making.

Most critiques of ProPublica’s COMPAS report followed the fair contest response, asserting that ProPublica focused on the wrong definition of fairness (Corbett-Davies et al., 2017; Dieterich et al., 2016; Flores et al., 2016; Gong, 2016). These respondents argued that COMPAS is fair because it satisfies sufficiency: each COMPAS score implies a similar likelihood of being arrested for both Black and white defendants. COMPAS produces a higher false positive rate for Black defendants simply because Black defendants are more likely to recidivate, not because COMPAS is racially biased. Most notably, Northpointe emphasized that the violation of separation presented by ProPublica “does *not* show evidence of bias, but rather is a natural consequence of using unbiased scoring rules for groups that happen to have different distributions of scores” (Dieterich et al., 2016).

The fair contest response seems appropriate within the lens of formal algorithmic fairness, which constrains analysis to the decision-making process alone. Given this scope of analysis, it seems fair to make decisions about people based on their likelihoods to exhibit a particular outcome of interest. For instance, if a Black and a white pretrial defendant are both high risk to be arrested in the future, then they should both be detained. Under this logic, algorithmic bias is a problem of systemic misrepresentation (e.g., over-predicting the risk of Black defendants relative to the

ground truth). It would be wrong for a high risk Black defendant to be detained while an equally high risk white defendant is released. Therefore, the best way to advance algorithmic fairness is to increase prediction accuracy and thereby ensure that decisions are based on accurate judgments about each individual (Hellman, 2020; Kleinberg et al., 2019).

However, because the fair contest response considers only the inputs and outputs of a specific algorithm, it fails to account for—and thus reproduces—broader patterns of injustice. First, the fair contest response fails to consider the presence of social hierarchies. In the case of risk assessments, Black and white defendants do not just “happen to have different distributions of scores,” as adherents of sufficiency assert (Dieterich et al., 2016). Instead, past and present discrimination has created social conditions in the US in which Black people are empirically at higher risk to commit crimes (Cooper & Smith, 2011; Sampson et al., 2005).<sup>9</sup> This disparity results from social oppression rather than from differences in inherent criminality (Muhammad, 2011). For instance, discriminatory practices such as segregation (Rothstein, 2017), racial criminalization (Butler, 2017; Muhammad, 2011), and severe underfunding of schools (EdBuild, 2019) all increase crime (Krivo et al., 2009; Lochner & Moretti, 2004; Rose & Clear, 1998).

Second, the fair contest response ignores the consequences of the actions that the algorithm informs. When a risk assessment labels a defendant “high risk,” that person is likely to be detained in jail until their trial. This practice of detaining defendants due to their crime risk, known as “preventative detention,” is both controversial and harmful. When the US Supreme Court deemed preventative detention constitutional in 1987, Justice Thurgood Marshall declared the practice “incompatible with the fundamental human rights protected by our Constitution” (U.S. Supreme Court, 1987). Preventative detention has faced continued scrutiny and challenge for undermining the rights of the accused and exacerbating mass incarceration (Baradaran, 2011; Koepke & Robinson, 2018). Pretrial detention imposes severe costs on defendants, including the loss of freedom, an increased likelihood of conviction, and a reduction in future employment (Dobbie et al., 2018).

By failing to account for the social hierarchies and unjust policies associated with pretrial decision-making, the fair contest response suggests a reform strategy in which even the best-case scenario—a perfectly accurate risk assessment—would perpetuate racial inequity.<sup>10</sup> Because Black defendants recidivate at higher rates than white defendants (Cooper & Smith, 2011; Flores et al., 2016; Larson et al., 2016; Sampson et al., 2005), a perfect risk assessment will accurately label a higher proportion of Black defendants as “high risk.” After all, if data is collected about an unequal

---

<sup>9</sup> This is true above and beyond racial disparities in arrest and enforcement patterns. Measurement bias is typically present in crime datasets but is not the only source of racial disparities in crime rates. For a broader discussion of the relationship between measurement and algorithmic fairness, see (Jacobs & Wallach, 2021).

<sup>10</sup> Because this risk assessment makes perfect predictions, it would satisfy both sufficiency and separation (Kleinberg et al., 2016).

society, then an accurate algorithm trained on that data will reflect those unequal conditions. To the extent that these predictions direct pretrial decisions, this risk assessment would lead to a higher pretrial detention rate for Black defendants than white defendants. This would, in effect, punish Black communities for having been unjustly subjected to criminogenic circumstances in the first place.

In sum, the fair contest response provides a meager guide for achieving algorithmic justice in settings with substantive inequalities. The central injustice of risk assessments is not that flawed data might lead an algorithm to make erroneous predictions of someone's crime risk. Instead, the central injustice is that racial stratification makes Black defendants higher risk than white defendants and that being deemed high risk leads to the loss of liberty. The fair contest response thus provides only the limited benefits of formal equality. Although a perfect risk assessment may help some Black defendants who are low risk but could be stereotyped as high risk, it would also naturalize the fact that many Black defendants actually are high risk and become incarcerated as a result.

#### **4.2 The Formalism Response: Constraining Reform**

The second formal algorithmic fairness response to the impossibility of fairness is what I call the “formalism response.” Recognizing that sufficiency reflects a limited notion of fairness, the formalism response does not require strict adherence to this measure. Instead, the formalism response focuses on analyzing the tradeoffs between notions of fairness. In particular, the formalism response suggests using the explicit mathematical formalization required by algorithms to rigorously consider the tradeoffs between separation and sufficiency in any given context.<sup>11</sup>

Under the formalism response, the formalism of algorithms provides a reality check by revealing the difficult tradeoffs between notions of fairness that might otherwise remain opaque and unarticulated (Barocas et al., 2019; Berk et al., 2018; Ligett, 2021). Algorithms provide “clarity” to help us identify and manage the unavoidable tradeoffs between competing goals (Kleinberg et al., 2019; Sunstein, 2019). Proponents of this view argue that algorithms can “be a positive force for social justice” because they “let us precisely *quantify tradeoffs* among society's different goals” and “force us to make more explicit judgments about underlying principles” (Kleinberg et al., 2019).

As with the fair contest response, the formalism response appears appropriate through the lens of formal algorithmic fairness, which constrains analysis to the decision-making process alone. Given this scope of analysis, it seems that reform interventions are limited to tuning the tradeoffs between

---

<sup>11</sup> The formalism response is inclusive of the fair contest response: after considering the tradeoffs, one could determine that an algorithm should be optimized for sufficiency. The formalism response can also account for other tradeoffs, such as the tension between fairness and accuracy.

values within the given decision-making procedure. Under this logic, the mathematical specificity of algorithms makes it possible to consider these tradeoffs more rigorously. For instance, given an existing population of Black and white defendants, reform involves grappling with the tradeoffs between sufficiency and separation in pretrial adjudication. Therefore, the best way to advance algorithmic fairness is to identify the most desirable balance between sufficiency and separation based on the particular context at hand.

However, because the formalism response limits analysis to individual decision-making processes, it yields an incredibly narrow range of possible reforms. First, the formalism response leaves us stuck making a zero-sum choice between two limited notions of fairness. Although separation may appear to be a desirable alternative to sufficiency, separation also fails to account for subordination. In the case of risk assessments, separation entails having different thresholds for Black and white defendants (e.g., a higher risk threshold for labeling Black defendants “high risk”). This practice would seem to obviate the point of using algorithmic risk predictions at all, as risk scores would have different meanings based on a defendant’s race (Flores et al., 2016; Mayson, 2019). Such explicit differential treatment based on race would be illegal to implement in many instances (Corbett-Davies et al., 2017; Hellman, 2020). Furthermore, although a lack of separation demonstrates that different groups face disparate burdens from mistaken judgments (Chouldechova, 2017; Hellman, 2020), separation does not prevent the injustices associated with accurate predictions. As demonstrated by the perfect pretrial risk assessment described in Section 4.1, an algorithm can satisfy separation while still reproducing racial hierarchy.

Second, the formalism response suggests a constrained and techno-centric reform strategy. Although the formalism response provides “clarity” regarding the tradeoffs involved in promoting fairness, this clarity is limited to the narrow scope of specific decision-making procedures. Everything beyond this scope is treated as static and thus irrelevant to evaluations of fairness. For instance, research on fairness in risk assessments explicitly places structural disadvantage and racial disparities outside the scope of algorithms and the responsibility of developers (Chouldechova, 2017; Corbett-Davies et al., 2017; Kleinberg et al., 2019). Following this logic, the formalism response suggests that implementing an algorithm is the only possible (or, at least, pertinent) alternative to the status quo (Berk et al., 2018; Kleinberg et al., 2019; Miller, 2018). This leads to the conclusion that the only appropriate path for reform is to improve specific decision-making processes using algorithms.

This approach is fundamentally limited as a strategy for achieving equitable public policy: egalitarian goals can rarely be achieved by reforming only the mechanisms of specific decision points. Reforms that aim to remedy structural oppression by targeting decision-making procedures often obscure and entrench the actual sources of oppression (Kahn, 2017; Murakawa, 2014). In the criminal justice system, for instance, “[a]dministrative tinkering does not confront the damning

features of the American carceral state, its scale and its racial concentration” (Murakawa, 2014). Implementing a pretrial risk assessment thus legitimizes preventative detention and hinders efforts to promote less carceral alternatives (Green, 2020).

In fact, the narrow purview of the formalism response is what makes the impossibility of fairness appear to be such an intractable dilemma. It is only because analysis is restricted to decision-making procedures that the tension between fairness definitions is interpreted as a fundamental “impossibility of fairness.” Mathematical proofs demonstrate that it is impossible to satisfy all mathematical definitions of fairness when making decisions about individuals in an unequal society. What is strictly “impossible” is simultaneously achieving two different mathematical notions of fair decision-making. However, by limiting analysis to isolated decision points, formal algorithmic fairness magnifies the stakes of this mathematical incompatibility, turning a constraint on fair decision-making into a constraint on fairness writ large. When all other aspects of society are treated as static or irrelevant, an algorithm’s behavior comes to represent “total fairness” (Berk et al., 2018). Under this assumption, the zero-sum tradeoff between mathematical definitions of fair decision-making represents an inescapable limitation on “total fairness.”

In sum, although the formalism response provides mathematical rigor about tradeoffs within particular decision points, it also obscures more systemic pathways for reform. The central question for reform is not simply how to tweak a particular decision-making process. Instead, the central question is how alter policies and institutions in ways that promote substantive equality. Because formal algorithmic fairness ignores reforms outside of specific decision points, it takes those reforms off the table at the outset of analysis. Although this strategy might be sufficient in some settings, it fails in the policy settings where concerns about injustice are most salient.

#### **4.3 Recap: The Methodological Limits of Formal Algorithmic Fairness**

Formal algorithmic fairness is methodologically incapable of promoting justice in policy settings with entrenched inequality. Akin to formal equality, formal algorithmic fairness formulates fairness within the scope of isolated decision points. As a result, formal algorithmic fairness is unable to account for social hierarchies and policies that exacerbate those hierarchies. Formal algorithmic fairness therefore traps algorithmic reform efforts within the impossibility of fairness, suggesting techno-centric reforms that entrench injustice. In Anderson’s terms, formal algorithmic fairness “can’t help us identify solutions” that address injustices because it fails to “represent the causes of [those] injustices” (Anderson, 2009). In Dewey’s terms, the issues with “what specific suggestions are entertained and which are dismissed” under formal algorithmic fairness are due to “[t]he way in which the problem is conceived” (Dewey, 1938). Thus, in order to develop a positive agenda for algorithmic justice, it is necessary to develop a new methodology for algorithmic fairness grounded in substantive equality.

## **5 SUBSTANTIVE ALGORITHMIC FAIRNESS: ESCAPING THE IMPOSSIBILITY OF FAIRNESS**

Given the methodological limits of formal algorithmic fairness, this section focuses on the second task of reforming algorithmic fairness: developing an alternative approach that operationalizes a social justice orientation into algorithmic fairness. In order to develop this method, I draw on substantive equality (described in Section 3.1) and the relational and structural reform strategies (described in Section 3.2).

As an alternative to formal algorithmic fairness, I propose a method of “substantive algorithmic fairness.” Substantive algorithmic fairness is an approach to algorithmic fairness in which the scope of analysis encompasses the social hierarchies and institutional structures that surround particular decision points. The goal is not to incorporate substantive equality into a formal mathematical model. This approach of “formalist incorporation” may yield some benefits, but would be subject to many of the same limits as formal algorithmic fairness (Green & Viljoen, 2020). As with fairness (Binns, 2018; Green & Hu, 2018; Jacobs & Wallach, 2021; Selbst et al., 2019), reducing substantive equality to mathematical definitions would narrow and distort the concept. Substantive algorithmic fairness therefore follows an approach of “algorithmic realism” (Green & Viljoen, 2020), incorporating algorithms into broader agendas for pursuing substantive equality. However, substantive algorithmic fairness does not entirely reject formal algorithmic fairness. Instead, it represents an expansion of algorithmic fairness methods, adopting substantive equality tools to reason about when formal algorithmic fairness is (and is not) appropriate.

Because of its broad frame of analysis, substantive algorithmic fairness provides a guide for using algorithms to promote equitable public policy without being constrained by the impossibility of fairness. Substantive algorithmic fairness reveals that the impossibility of fairness is a misnomer: when an algorithm confronts the incompatibility between fairness definitions, it suggests not that fairness is impossible writ large, but that algorithms are being used to pursue a misguided reform strategy. Debates and consternation about the impossibility of fairness are most extreme when making decisions in which a) an oppressed group disproportionately exhibits the attributes deemed “negative” in the given context (e.g., indicators of high crime risk), and b) policy punishes (or restricts benefits to) individuals who exhibit these negative attributes. When these relational and structural factors are present, any attempt to improve decision-making with an algorithm will confront the impossibility of fairness. The proper response to the impossibility of fairness is not to tinker within the contours of this intractable dilemma, but to reform the relational and structural factors that produce the dilemma. If there were no social hierarchies or if consequential decisions did not exacerbate social hierarchies, then the impossibility of fairness would not arise (or, at least, would not be so concerning). In other words, substantive algorithmic fairness indicates that reforms should target relational and structural inequalities, not just the precise mechanisms of decision-making.

This section proceeds in three parts. First, I describe the general principles of substantive algorithmic fairness. Second, I apply substantive algorithmic fairness to pretrial reform. Third, I describe what substantive algorithmic fairness entails as a practice for algorithm developers.

### **5.1 The Substantive Algorithmic Fairness Approach to Reform**

As with formal algorithmic fairness, the starting point for reform in substantive algorithmic fairness is concern about discrimination or inequality within a particular decision-making process. Drawing on the substantive equality approaches introduced in Section 3, substantive algorithmic fairness presents a three-step strategy for promoting equality in such scenarios. Each step can be boiled down to a central question. 1) What is the substance of the inequalities in question? 2) What types of reforms can remediate the identified substantive inequalities? 3) What roles, if any, can algorithms play to enhance or facilitate the identified reforms?

The first step is to diagnose the substance of the inequalities in question. This entails looking for conditions of hierarchy and questioning how social and institutional arrangements reinforce those conditions (MacKinnon, 2011). When faced with disparities in data, substantive algorithmic fairness asks: do these disparities reflect social conditions of hierarchy? Similarly, when faced with particular decision points, substantive algorithmic fairness asks: do these decisions (and the interventions that they facilitate) exacerbate social hierarchies? If the answers to both questions are no, then formal algorithmic fairness presents an appropriate path forward. However, if the answers to these questions are yes—as they often will be when confronting inequalities in high-stakes decisions—then reforms limited to decision-making processes will be insufficient.

The second step is to consider what types of reforms can remediate the identified substantive inequalities. Substantive algorithmic fairness draws on the reforms proposed by Anderson (Anderson, 1999), Minow (Minow, 1991), and Fishkin (Fishkin, 2014) for promoting equality without becoming trapped by intractable dilemmas. The first approach is the relational response: reform the relationships that create and sustain social hierarchies. The second approach is the structural response: reshape the structure of decisions to avoid or lower the stakes of decisions that exacerbate social hierarchies. Because these reforms target the relational and structural factors that produce equality dilemmas, they are not subject to the impossibility of fairness.

The third step is to analyze whether and how algorithms can enhance or facilitate the reforms identified in the second step. The critical words here are “enhance” and “facilitate.” Rather than treating algorithms as the central component of reform, the analysis here should consider whether and how algorithms can support larger agendas for reform. Thus, in considering the potential role for algorithms, computer scientists should be wary of technological determinism and the assumption that algorithms can remedy all social problems. Algorithmic interventions should be

considered through an “agnostic approach” that prioritizes the reform agenda identified in the second step, without assuming any necessary or particular role for algorithms (Green & Viljoen, 2020). This approach requires decentering technology when studying injustice and remaining attentive to the broader structural forces of marginalization (Gangadharan & Niklas, 2019). In some cases, this will mean recognizing that algorithms are unnecessary or even detrimental tools for reform. In other cases, this will mean identifying new approaches for developing and applying algorithms to help combat oppression. Algorithms can play productive roles in support of broader efforts for social change (Abebe et al., 2020), particularly when deployed in conjunction with policy and governance reforms (Green, 2019).

## **5.2 Example: The Substantive Algorithmic Fairness Approach to Pretrial Reform**

We can see the benefits of substantive algorithmic fairness by considering how it applies in the context of pretrial reform. Formal algorithmic fairness suggests that the appropriate pretrial reform strategy is to make release/detain decisions using algorithmic predictions of risk. Despite the support for pretrial risk assessments among many engineers and policymakers, this approach upholds racial injustice and leaves decision-making caught within the impossibility of fairness. In contrast, substantive algorithmic fairness suggests reforms that more robustly challenge the injustices associated with pretrial decision-making and that provide an escape from the impossibility of fairness. Although this approach highlights the limits of pretrial risk assessments, it also suggests new paths for reform and new roles for algorithms.

When pursuing pretrial reform through substantive algorithmic fairness, the first step is to consider the substance of inequalities that manifest in pretrial decision-making. As described in Section 4.1, the disparity in recidivism rates across Black and white defendants reflects conditions of racial hierarchy. This disparity cannot be attributed to chance or to inherent group differences (nor is it solely the result of measurement bias). Furthermore, preventative detention exacerbates this hierarchy by depriving high-risk defendants of rights and subjecting them to a range of negative outcomes.

The second step is to consider what reforms could appropriately address the substantive inequalities identified in the first step. Here, we can follow the relational and structural responses. The relational response suggests altering the relationships that define “risk” and shape its unequal distribution across the population. This provides a clear contrast to pretrial risk assessments, which treat risk as an intrinsic and neutral attribute of individuals and thereby naturalize group differences in risk that are the product of oppression. The relational response provides an alternative approach: reform the social arrangements that make risk a socially salient category. The relational response thus suggests aiming to reduce the crime risk of Black communities by alleviating criminogenic conditions of disadvantage. For instance, public policies that extend access to education (Lochner & Moretti, 2004), welfare (Tuttle, 2019), and affordable housing (Diamond & McQuade, 2019)

all reduce crime, and therefore could reduce the racial disparity in crime risk. The relational response also suggests combatting the association of Blackness with criminality and the effects of this association. This entails not merely challenging stereotypes that link Blackness with crime, but also decriminalizing behaviors that were previously criminalized to subjugate minorities (Butler, 2017; Muhammad, 2011).

The structural response suggests altering the structure of decisions to reduce the harmful consequences associated with being high risk to recidivate. This provides a clear contrast to pretrial risk assessments, which uphold the notion that the appropriate response to high-risk defendants is incarceration. The structural response provides an alternative approach: reform policy to ensure that being high risk no longer prompts such severe punishment. The structural response thus suggests attempting to minimize the scope and harms of decisions that determine one's freedom and opportunities based on their risk of recidivism. If fewer people were subjected to decisions in which liberty and well-being depend on exhibiting low levels of crime risk, racial disparities in the distribution of risk would be less consequential. Most directly, such an approach could entail abolishing (or drastically reducing the scope of) pretrial detention, such that fewer people would be incarcerated, regardless of their risk level. Reforms could also aim to decrease the downstream damages of pretrial detention. For instance, reducing the effects of pretrial detention on increased conviction and diminished future employment would reduce the harms associated with being high risk. Another reform along these lines would be to shift from responding to risk with punishment to responding with social or material support, such that the consequence of being high risk is to receive aid rather than incarceration.

The third step is to consider the potential role for algorithms in advancing relational and structural reforms. In some cases, this analysis will provide arguments against the use of certain algorithms for reform. For instance, because pretrial risk assessments naturalize racial disparities in risk that are the product of oppression and legitimize preventative detention, these algorithms conflict with the relational and structural responses. In other cases, however, this analysis will reveal new, fruitful roles for algorithms in pretrial reform. Importantly, however, these alternative roles for algorithms will involve a broader scope than just the pretrial decision-making process and will operate in conjunction with other reforms.

Following the relational response, the key question is whether algorithms can enhance or facilitate the identified relational reforms. One path along these lines involves using algorithms to reduce the crime risk of Black communities by alleviating criminogenic conditions of disadvantage. For instance, algorithms have been used to increase access to education (Lakkaraju et al., 2015), welfare (DataSF, 2018), and affordable housing (Ye et al., 2019), all of which can reduce the crime risk of disadvantaged groups. Another direction involves using algorithms to combat the criminalization of minorities. Several states have implemented algorithms to streamline the

process of expunging criminal records, which is likely to disproportionately benefit minority and low-income individuals (Johnston, 2022). Similarly, statistical analyses have helped to document how stop-and-frisk discriminates against minorities and to push for altering or abolishing this practice (Denvir, 2015; Goel et al., 2016).

Following the structural response, the key question is whether algorithms can enhance or facilitate the identified structural reforms. One path along these lines involves using algorithms to reduce the harms of the racial disparity in recidivism risk. Algorithms can be used to target supportive rather than punitive responses to risk (Barabas et al., 2018; Mayson, 2019), thus mitigating rather than compounding the injustices behind the high recidivism risk of Black defendants. Another direction involves using algorithms to support broader political agendas for structural reforms. For instance, algorithms could help justify structural reforms by exposing the false promises of pretrial risk assessments (Angwin et al., 2016; Green & Chen, 2019) and by providing a systemic view of how the criminal justice system exacerbates racial inequalities (Crespo, 2015). Algorithms could also be used to make structural reforms more possible by empowering communities advocating for criminal justice reform and supporting the campaigns of political candidates promising such reforms.

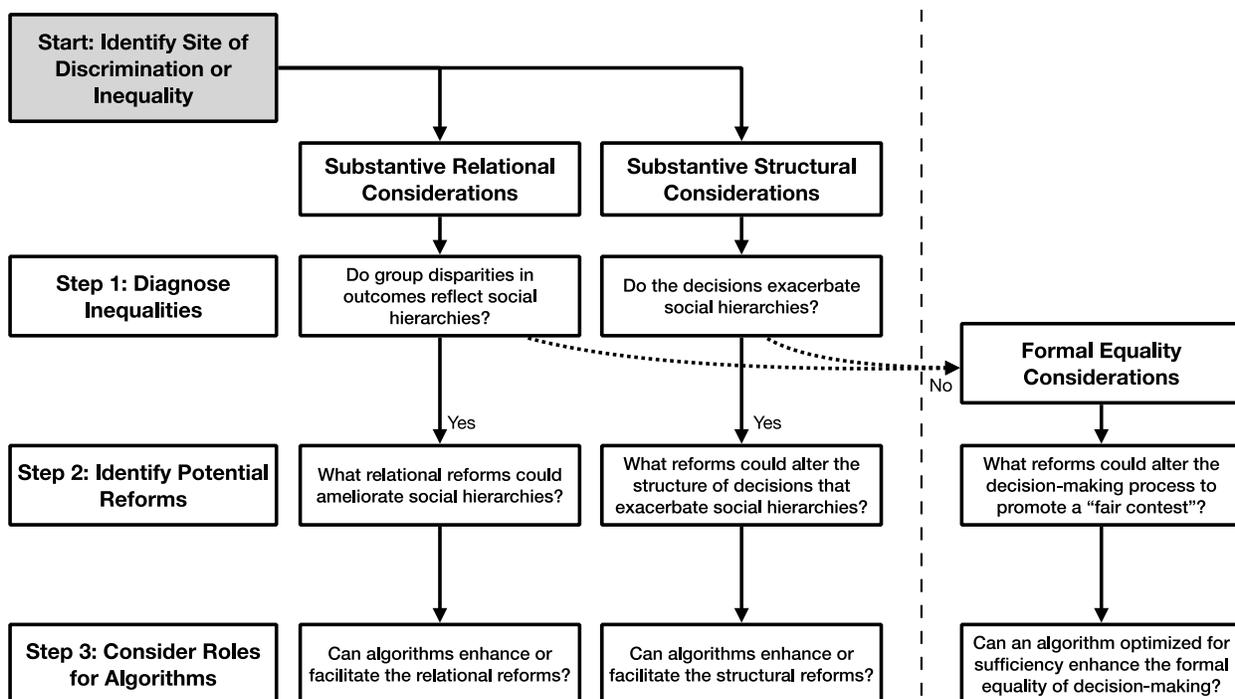
In sum, substantive algorithmic fairness demonstrates how an expansive analysis of social conditions and institutions can lead to rigorous theories of social change, and how those theories of change can inform work on algorithms that is not subject to the impossibility of fairness. Starting with these broader reform agendas provides paths for algorithms in pretrial reform that involve more than just pretrial risk assessments. It is important to note that none of these alternative algorithmic interventions would completely solve the problems of pretrial detention—that is an unrealistic goal for any individual reform. Nor are algorithms necessarily the centerpiece of reform. Instead, these algorithmic interventions operate in conjunction with other reforms, aiming to enhance efforts pushing for substantive pretrial reform. These benefits could accrue similarly in other areas in which the impossibility of fairness has been interpreted as a significant and intractable barrier to algorithmic justice, such as child welfare (Chouldechova et al., 2018) and college admissions (Friedler et al., 2021).

### **5.3 Substantive Algorithmic Fairness in Practice**

Substantive algorithmic fairness offers a new direction for algorithmic fairness. It shifts the field's concern away from formal mathematical models of “fair” decision-making and toward substantive evaluations of how algorithms can (and cannot) combat social hierarchies. In doing so, substantive algorithmic fairness aligns the field with recent calls for algorithmic “justice” (Bui & Noble, 2020; Costanza-Chock, 2020; Green, 2018), “equity” (D’Ignazio & Klein, 2020), and “reparation” (Davis et al., 2021).

Substantive algorithmic fairness provides a new guide for computer scientists hoping to promote equity with algorithms. When considering injustices that they would like to help ameliorate, computer scientists can follow substantive algorithmic fairness by working through the flowchart in Figure 1. This flowchart translates substantive equality goals into concrete questions for computer scientists to consider before developing an algorithm. In turn, the flowchart can direct computer scientists away from the narrow and techno-centric reforms typically suggested by formal algorithmic fairness and related efforts to pursue “social good” (Green, 2021). It can also help computer scientists identify when narrower, formal algorithmic fairness methods may actually be appropriate.

**Figure 1. Flowchart to support computer scientists in following substantive algorithmic fairness in practice.** The process begins at the top of the flowchart, with concern about discrimination or inequality in a particular decision-making process. This feeds into the substantive equality considerations focused on relational and structural inequalities. If neither relational nor structural concerns are salient (i.e., the answers to both questions in Step 1 are “No”), then the process transitions to formal equality considerations. In this case, the tasks for computer scientists resemble those that already exist within formal algorithmic fairness. In this sense, substantive algorithmic fairness represents an expansion of algorithmic fairness methodology rather than a complete rejection of formal algorithmic fairness.



Of course, answering the flowchart's questions can be a difficult and politically contested task. The flowchart's questions lack straightforward, objective answers and should not be made by computer scientists alone. As a result, substantive algorithmic fairness requires new practices and training for computer scientists hoping to improve public policy. The questions in Step 1 require engagement with philosophical and social scientific literature about inequality. In order to answer these questions, computer scientists must gain rigorous knowledge about the social and policy context in question. These efforts will be significantly enhanced by interdisciplinary collaborations that bring scholars from law, sociology, philosophy, and other fields into the diagnoses of inequality. It is also essential to engage with communities affected by the discrimination or inequality that motivates a computer scientist's concerns.

The questions in Step 2 require engagement with scholarship on social change as well as communities advocating for reform. As in Step 1, interdisciplinary collaborations will help computer scientists identify reform strategies that could be effective and achievable in the given context. It is also particularly important for computer scientists to engage with communities already advocating for reform. Substantive algorithmic fairness requires broad deliberation that includes the communities most likely to be affected by any reform or algorithm. Prior work has demonstrated that the choice of fairness metric is a political task that should be made democratically (Wong, 2020). Substantive algorithmic fairness suggests that democratizing algorithmic fairness requires an even broader scope. In addition to democratizing the choice of technical parameters, it is necessary to democratize decisions such as which reforms to pursue and whether to consider using algorithms at all.

The questions in Step 3 require engagement with scholars of sociotechnical systems and communities advocating for reform. Research from fields such as STS (science, technology, and society) can help computer scientists understand how technology affects society and how efforts to improve society with technology can go wrong. As above, it is also necessary to engage with communities to understand whether a potential algorithm aligns with their reform agendas. This is essential given the tendency of computer scientists to see technology as a solution to social issues. The primary goal in Step 3 is not to develop an algorithm that achieves a particular reform on its own. Instead, the goal is to develop an algorithm that can advance broader reform agendas. Recent work provides several examples of how data analysis and technology design can be incorporated into community-driven reform efforts that challenge oppression (Asad, 2019; Costanza-Chock, 2020; Lewis et al., 2018; Maharawal & McElroy, 2018; Meng & DiSalvo, 2018).

As these three steps suggest, substantive algorithmic fairness requires new types of training for computer scientists. Computer science training must expand beyond the traditional emphasis on the mathematical properties of algorithms to incorporate normative reasoning, sociotechnical systems, and theories of social change. It is also important that computer science training inculcate

a focus on the real-world social impacts of algorithms. In addition to courses focused on ethics and sociotechnical systems, curricula should incorporate practice-based classes in which students collaborate with organizations (e.g., government agencies, nonprofits, and advocacy organizations). Such courses can help students consider an algorithm's impacts in light of broader social contexts and appreciate the power of choosing to not design systems that could exacerbate inequality (Graeff, 2020).

Through these questions and practices, substantive algorithmic fairness provides a compass to help computer scientists reason about the appropriate roles for algorithms in combatting inequity. Debates about algorithmic reforms often feature a binary contest between algorithmic reforms and the status quo: when critics challenge the use of algorithms, proponents for algorithms argue that the only alternative to implementing fallible and biased algorithms is to rely on even more fallible and biased human decision-makers (Berk et al., 2018; Kleinberg et al., 2019; Miller, 2018). Substantive algorithmic fairness demonstrates that reformers need not accept this binary choice between implementing a superficially "fair" algorithm and leaving the status quo in place. Instead, substantive algorithmic fairness presents strategies for integrating algorithms into broader efforts to pursue equity. And although substantive algorithmic fairness begins with an ambitious (some might say utopian) vision of substantive equality, the reforms that it suggests are all incremental: none will create a substantively equal society on their own. Each reform, however, moves society one step closer to substantive equality. In this sense, substantive algorithmic fairness takes after political theories of "non-reformist reforms" (Gorz, 1967), "real utopias" (Wright, 2010), and prison abolition (McLeod, 2015), all of which present strategies for linking short-term, piecemeal reforms with long-term, radical agendas for social justice.

Despite these benefits, substantive algorithmic fairness does not provide a precise roadmap for reform. It presents a sequence of questions, with conceptual tools for answering those questions in a principled manner, rather than a mandatory checklist. It cannot be reduced to an optimization problem. This lack of explicit prescription is not so much a limit of substantive algorithmic fairness as an inescapable reality of pursuing substantive social and political reform. There is no single or straightforward path for how to achieve change (Unger, 2005; Wright, 2010). The hardest political questions often revolve around which reforms to pursue in any specific situation, among many potential paths forward. Making these judgments requires contextual assessments of feasibility and impact as well as engagement with affected communities. In some settings, particularly where substantive concerns about social hierarchy and unjust policies are less severe, this analysis may even suggest reforms that align with formal algorithmic fairness. There similarly is no straightforward mechanism for determining how to best incorporate algorithms into reform efforts. Future work is necessary to better understand the appropriate roles for algorithms in reform efforts, the conditions that facilitate effective algorithmic reforms, and how to allocate authority over algorithmic reforms.

Furthermore, efforts to achieve substantive algorithmic fairness in practice face a variety of barriers. Many political actors and technology companies benefit from and promote formal algorithmic fairness, as it allows them to embrace “fairness” without making significant political or economic concessions (Bui & Noble, 2020; Green, 2020; Powles & Nissenbaum, 2018). Efforts to achieve the reforms suggested by substantive algorithmic fairness will often confront these forces opposed to structural change. The exclusion of women and minorities from algorithm development also leads to notions of algorithmic fairness that are inattentive to the lived realities of oppressed groups (West, 2020). Additionally, institutional barriers and incentives hinder the necessary types of interdisciplinary research and training. Thus, as with all efforts to achieve substantive equality, substantive algorithmic fairness requires ongoing political struggle to achieve conditions amenable to reform.

## **6 CONCLUSION**

Algorithmic fairness provides an increasingly influential toolkit for promoting equitable public policy. It is therefore essential to consider whether algorithmic fairness provides suitable conceptual and practical tools to guide reform. If algorithmic fairness methodology cannot comprehensively recognize and represent the nature of injustices, it will fail to identify effective paths for remediating those injustices. The impossibility fairness suggests that algorithmic fairness suffers from methodological deficits. In light of this concern, this article took on two essential tasks for reforming algorithmic fairness.

First, I diagnosed the current methodological limits of algorithmic fairness. The current methodology—formal algorithmic fairness—is poorly equipped for enhancing equality. Because it restricts analysis to isolated decision points, formal algorithmic fairness cannot account for social hierarchies and the impacts of decisions informed by algorithms. As a result, formal algorithmic fairness traps reform efforts within the impossibility of fairness and suggests reforms that uphold social hierarchies. Before algorithmic fairness can productively guide efforts to pursue equality, we must alter its methodology to encompass more comprehensive conceptual and practical tools.

Second, I proposed an alternative methodology for algorithmic fairness that overcomes the limits of formal algorithmic fairness. Substantive algorithmic fairness provides a new orientation for algorithmic fairness, incorporating algorithms into broader movements for reform. In doing so, substantive algorithmic fairness offers an escape from the impossibility of fairness and suggests new roles for algorithms in combatting oppression. In shifting away from formal mathematical models (and associated interventions such as pretrial risk assessments), this reorientation prompts a new positive agenda for how to act on recent calls to shift the field’s emphasis from “fairness” to “justice” (Bui & Noble, 2020; Costanza-Chock, 2020; Green, 2018), “equity” (D’Ignazio & Klein, 2020), and “reparation” (Davis et al., 2021).

Although substantive algorithmic fairness does not yield a precise roadmap for reform, it presents concrete steps to help computer scientists link visions of substantive equality with incremental algorithmic reforms. Substantive algorithmic fairness reveals that reform-minded computer scientists do not face a binary choice between implementing a “fair” algorithm and doing nothing. Instead, there are many potential reforms to consider—all of them, in some form, incremental—and many potential roles for algorithms to enable or supplement those reforms. Substantive algorithmic fairness provides a method to diagnose the inequalities in need of reform, evaluate which reforms can best advance substantive equality, and consider how algorithms can support those reforms.

No single reform—algorithmic or otherwise—can create a substantively equal society. However, algorithmic fairness researchers need not restrict themselves to a formal algorithmic fairness methodology that constrains opportunities for reform and often reinforces oppression. By starting from substantive accounts of social hierarchy and social change, the field of algorithmic fairness can stitch together incremental algorithmic reforms that collectively build a more egalitarian society.

### **Acknowledgments**

I am grateful to Elettra Bietti, Matt Bui, Ben Fish, Evan Green, Will Holub-Moorman, Lily Hu, Abbie Jacobs, Andrew Schrock, Salomé Viljoen, and Zach Wehrwein for valuable suggestions about how to improve this manuscript.

## **7 REFERENCES**

- Abebe, R., Barocas, S., Kleinberg, J., Levy, K., Raghavan, M., & Robinson, D. G. (2020). *Roles for computing in social change* Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, Barcelona, Spain. <https://doi.org/10.1145/3351095.3372871>
- Anderson, E. (2009). Toward a Non-Ideal, Relational Methodology for Political Philosophy: Comments on Schwartzman's *Challenging Liberalism*. *Hypatia*, 24(4), 130-145. [www.jstor.org/stable/20618184](http://www.jstor.org/stable/20618184)
- Anderson, E. S. (1999). What is the Point of Equality? *Ethics*, 109(2), 287-337.
- Angwin, J., & Larson, J. (2016). Bias in Criminal Risk Scores Is Mathematically Inevitable, Researchers Say. *ProPublica*. <https://www.propublica.org/article/bias-in-criminal-risk-scores-is-mathematically-inevitable-researchers-say>
- Angwin, J., Larson, J., Mattu, S., & Kirchner, L. (2016). Machine Bias. *ProPublica*. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- Arneson, R. (2013). Egalitarianism. *The Stanford Encyclopedia of Philosophy*. <https://plato.stanford.edu/entries/egalitarianism/>

- Arnold Ventures. (2019). Statement of Principles on Pretrial Justice and Use of Pretrial Risk Assessment. <https://craftmediabucket.s3.amazonaws.com/uploads/Arnold-Ventures-Statement-of-Principles-on-Pretrial-Justice.pdf>
- Asad, M. (2019). Prefigurative Design as a Method for Research Justice. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW). <https://doi.org/10.1145/3359302>
- Barabas, C., Virza, M., Dinakar, K., Ito, J., & Zittrain, J. (2018). Interventions over Predictions: Reframing the Ethical Debate for Actuarial Risk Assessment. Proceedings of the 1st Conference on Fairness, Accountability and Transparency, Proceedings of Machine Learning Research.
- Baradaran, S. (2011). Restoring the Presumption of Innocence. *Ohio State Law Journal*, 72, 723-776.
- Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness and Machine Learning*. fairmlbook.org.
- Berk, R., Heidari, H., Jabbari, S., Kearns, M., & Roth, A. (2018). Fairness in Criminal Justice Risk Assessments: The State of the Art. *Sociological Methods & Research*, 1-42. <https://journals.sagepub.com/doi/abs/10.1177/0049124118782533>
- Binns, R. (2018). *Fairness in Machine Learning: Lessons from Political Philosophy* Proceedings of the 1st Conference on Fairness, Accountability and Transparency, <http://proceedings.mlr.press>
- Bui, M. L., & Noble, S. U. (2020). We're Missing a Moral Framework of Justice in Artificial Intelligence: On the Limits, Failings, and Ethics of Fairness. In M. D. Dubber, F. Pasquale, & S. Das (Eds.), *The Oxford Handbook of Ethics of AI*. Oxford University Press.
- Butler, P. (2017). *Chokehold: Policing Black Men*. The New Press.
- California Legislature. (2021). AB-13 Public contracts: automated decision systems. [https://leginfo.legislature.ca.gov/faces/billTextClient.xhtml?bill\\_id=202120220AB13](https://leginfo.legislature.ca.gov/faces/billTextClient.xhtml?bill_id=202120220AB13)
- Chouldechova, A. (2017). Fair Prediction with Disparate Impact: A Study of Bias in Recidivism Prediction Instruments. *Big Data*, 5(2), 153-163.
- Chouldechova, A., Benavides-Prado, D., Fialko, O., & Vaithianathan, R. (2018). *A case study of algorithm-assisted decision making in child maltreatment hotline screening decisions* Proceedings of the 1st Conference on Fairness, Accountability and Transparency, Proceedings of Machine Learning Research. <https://proceedings.mlr.press/v81/chouldechova18a.html>
- Cooper, A., & Smith, E. L. (2011). Homicide Trends in the United States, 1980-2008. *U.S. Department of Justice, Bureau of Justice Statistics*. <https://www.bjs.gov/content/pub/pdf/htus8008.pdf>
- Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., & Huq, A. (2017). Algorithmic Decision Making and the Cost of Fairness. Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada.
- Costanza-Chock, S. (2020). *Design Justice: Community-Led Practices to Build the Worlds We Need*. MIT Press.

- Crenshaw, K. W. (1988). Race, Reform, and Retrenchment: Transformation and Legitimation in Antidiscrimination Law. *Harvard Law Review*, 101(7), 1331-1387.
- Crespo, A. M. (2015). Systemic Facts: Toward Institutional Awareness in Criminal Courts. *Harvard Law Review*, 129, 2049-2117.
- D'Ignazio, C., & Klein, L. F. (2020). *Data Feminism*. MIT Press.
- DataSF. (2018). Keeping Moms and Babies in Nutrition Program. <https://datasf.org/showcase/datascience/keeping-moms-and-babies-in-nutrition-program/>
- Davis, J. L., Williams, A., & Yang, M. W. (2021). Algorithmic reparation. *Big Data & Society*, 8(2). <https://doi.org/10.1177/20539517211044808>
- Denvir, D. (2015). The Key Ingredient in Stop-and-Frisk Reform: Open Data. *Bloomberg*. <https://www.bloomberg.com/news/articles/2015-08-24/the-missing-ingredient-in-stop-and-frisk-reform-open-data>
- Dewey, J. (1938). *Logic: The Theory of Inquiry*. Henry Holt and Company.
- Diamond, R., & McQuade, T. (2019). Who Wants Affordable Housing in Their Backyard? An Equilibrium Analysis of Low-Income Property Development. *Journal of Political Economy*, 127(3), 1063-1117. <https://doi.org/10.1086/701354>
- Dieterich, W., Mendoza, C., & Brennan, T. (2016). COMPAS Risk Scales: Demonstrating Accuracy Equity and Predictive Parity. *Northpointe Inc. Research Department*. [https://go.volarisgroup.com/rs/430-MBX-989/images/ProPublica\\_Commentary\\_Final\\_070616.pdf](https://go.volarisgroup.com/rs/430-MBX-989/images/ProPublica_Commentary_Final_070616.pdf)
- Dobbie, W., Goldin, J., & Yang, C. S. (2018). The Effects of Pretrial Detention on Conviction, Future Crime, and Employment: Evidence from Randomly Assigned Judges. *American Economic Review*, 108(2), 201-240. <https://doi.org/10.1257/aer.20161503>
- Doctorow, C. (2016). Algorithmic risk-assessment: hiding racism behind "empirical" black boxes. *Boing Boing*. <https://boingboing.net/2016/05/24/algorithmic-risk-assessment-h.html>
- EdBuild. (2019). \$23 Billion. <https://edbuild.org/content/23-billion/full-report.pdf>
- Eubanks, V. (2018). *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor*. St. Martin's Press.
- European Commission. (2021). Proposal for a Regulation laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). <https://digital-strategy.ec.europa.eu/en/library/proposal-regulation-laying-down-harmonised-rules-artificial-intelligence-artificial-intelligence>
- Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., & Venkatasubramanian, S. (2015). Certifying and Removing Disparate Impact. 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia.
- Fishkin, J. (2014). *Bottlenecks: A New Theory of Equal Opportunity*. Oxford University Press.
- Flores, A. W., Bechtel, K., & Lowenkamp, C. T. (2016). False Positives, False Negatives, and False Analyses: A Rejoinder to "Machine Bias: There's Software Used Across the Country

- to Predict Future Criminals. And It's Biased Against Blacks.”. *Federal Probation*, 80, 38-46.
- Friedler, S. A., Scheidegger, C., & Venkatasubramanian, S. (2021). The (Im)possibility of Fairness: Different Value Systems Require Different Mechanisms For Fair Decision Making. *Communications of the ACM*, 64(4), 136–143. <https://doi.org/10.1145/3433949>
- Gangadharan, S. P., & Niklas, J. (2019). Decentering technology in discourse on discrimination. *Information, Communication & Society*, 22(7), 882-899. <https://doi.org/10.1080/1369118X.2019.1593484>
- Goel, S., Rao, J. M., & Shroff, R. (2016). Precinct or prejudice? Understanding racial disparities in New York City's stop-and-frisk policy. *The Annals of Applied Statistics*, 10(1), 365-394.
- Gong, A. (2016). Ethics for powerful algorithms (1 of 4). *Medium*. <https://medium.com/@AbeGong/ethics-for-powerful-algorithms-1-of-3-a060054efd84>
- Gorz, A. (1967). *Strategy for Labor*. Beacon Press.
- Gosepath, S. (2021). Equality. *The Stanford Encyclopedia of Philosophy*. <https://plato.stanford.edu/entries/equality/>
- Government of Canada. (2021). Directive on Automated Decision-Making. <https://www.tbs-sct.gc.ca/pol/doc-eng.aspx?id=32592>
- Graeff, E. (2020). The Responsibility to Not Design and the Need for Citizen Professionalism. *Tech Otherwise*. <https://techotherwise.pubpub.org/pub/vizamy14>
- Green, B. (2018). Putting the J(ustice) in FAT. *Berkman Klein Center Collection - Medium*. <https://medium.com/berkman-klein-center/putting-the-j-ustice-in-fat-28da2b8eae6d>
- Green, B. (2019). *The Smart Enough City: Putting Technology in Its Place to Reclaim Our Urban Future*. MIT Press.
- Green, B. (2020). The False Promise of Risk Assessments: Epistemic Reform and the Limits of Fairness. Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency,
- Green, B. (2021). Data Science as Political Action: Grounding Data Science in a Politics of Justice. *Journal of Social Computing*.
- Green, B., & Chen, Y. (2019). Disparate Interactions: An Algorithm-in-the-Loop Analysis of Fairness in Risk Assessments. Proceedings of the Conference on Fairness, Accountability, and Transparency,
- Green, B., & Hu, L. (2018). The Myth in the Methodology: Towards a Recontextualization of Fairness in Machine Learning. Machine Learning: The Debates workshop at the 35th International Conference on Machine Learning,
- Green, B., & Viljoen, S. (2020). Algorithmic Realism: Expanding the Boundaries of Algorithmic Thought. Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency,
- Hardt, M., Price, E., & Srebro, N. (2016). Equality of Opportunity in Supervised Learning. 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

- Harris, K., & Paul, R. (2017). Pretrial Integrity and Safety Act of 2017. *115th Congress*.
- Hellman, D. (2020). Measuring Algorithmic Fairness. *Virginia Law Review*, *106*(4), 811-866.
- Hoffmann, A. L. (2019, 2019/06/07). Where fairness fails: data, algorithms, and the limits of antidiscrimination discourse. *Information, Communication & Society*, *22*(7), 900-915. <https://doi.org/10.1080/1369118X.2019.1573912>
- Jacobs, A. Z., & Wallach, H. (2021). Measurement and Fairness. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 375–385. <https://doi.org/10.1145/3442188.3445901>
- Johnston, R. (2022). New algorithm to expunge 500,000 criminal records in Utah. *StateScoop*. <https://statescoop.com/utah-record-expunge-algorithm-code-for-america/>
- Kahn, J. (2017). *Race on the Brain: What Implicit Bias Gets Wrong about the Struggle for Racial Justice*. Columbia University Press.
- Kalluri, P. (2020). Don't ask if artificial intelligence is good or fair, ask how it shifts power. *Nature*, *583*, 169. <https://doi.org/10.1038/d41586-020-02003-2>
- Kleinberg, J., Ludwig, J., Mullainathan, S., & Sunstein, C. R. (2019). Discrimination in the Age of Algorithms. *Journal of Legal Analysis*, *10*, 113-174. <https://doi.org/10.1093/jla/laz001>
- Kleinberg, J., Mullainathan, S., & Raghavan, M. (2016). Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*.
- Koepke, J. L., & Robinson, D. G. (2018). Danger Ahead: Risk Assessment and the Future of Bail Reform. *Washington Law Review*, *93*, 1725-1807.
- Krivo, L. J., Peterson, R. D., & Kuhl, D. C. (2009). Segregation, Racial Structure, and Neighborhood Violent Crime. *American Journal of Sociology*, *114*(6), 1765-1802. <https://doi.org/10.1086/597285>
- Lakkaraju, H., Aguiar, E., Shan, C., Miller, D., Bhanpuri, N., Ghani, R., et al. (2015). *A Machine Learning Framework to Identify Students at Risk of Adverse Academic Outcomes* Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia. <https://doi.org/10.1145/2783258.2788620>
- Larson, J., Mattu, S., Kirchner, L., & Angwin, J. (2016). How We Analyzed the COMPAS Recidivism Algorithm. *ProPublica*. <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
- Lewis, T., Gangadharan, S. P., Saba, M., & Petty, T. (2018). Digital Defense Playbook: Community Power Tools for Reclaiming Data. *Our Data Bodies*. [https://www.odbproject.org/wp-content/uploads/2019/03/ODB\\_DDP\\_HighRes\\_Single.pdf](https://www.odbproject.org/wp-content/uploads/2019/03/ODB_DDP_HighRes_Single.pdf)
- Ligett, K. (2021). FAccT 2021 Keynote: In Praise of Flawed Mathematical Models. <https://www.youtube.com/watch?v=gZrZwF3XDBw>
- Lochner, L., & Moretti, E. (2004). The Effect of Education on Crime: Evidence from Prison Inmates, Arrests, and Self-Reports. *American Economic Review*, *94*(1), 155-189.

- MacKinnon, C. A. (2011). Substantive Equality: A Perspective. *Minnesota Law Review*, 96, 1.
- MacKinnon, C. A. (2016). Substantive equality revisited: A reply to Sandra Fredman. *International Journal of Constitutional Law*, 14(3), 739-746.
- Maharawal, M. M., & McElroy, E. (2018, 2018/03/04). The Anti-Eviction Mapping Project: Counter Mapping and Oral History toward Bay Area Housing Justice. *Annals of the American Association of Geographers*, 108(2), 380-389. <https://doi.org/10.1080/24694452.2017.1365583>
- Mayson, S. G. (2019). Bias In, Bias Out. *Yale Law Journal*, 128(8), 2218-2300.
- McLeod, A. M. (2015). Prison Abolition and Grounded Justice. *UCLA Law Review*, 62, 1156-1239.
- Meng, A., & DiSalvo, C. (2018). Grassroots resource mobilization through counter-data action. *Big Data & Society*, 5(2), 2053951718796862. <https://doi.org/10.1177/2053951718796862>
- Merriam-Webster. (2021). Methodology. <https://www.merriam-webster.com/dictionary/methodology>
- Miller, A. P. (2018). Want Less-Biased Decisions? Use Algorithms. *Harvard Business Review*. <https://hbr.org/2018/07/want-less-biased-decisions-use-algorithms>
- Minow, M. (1991). *Making All the Difference: Inclusion, Exclusion, and American Law*. Cornell University Press.
- Minow, M. (2021). Equality vs. Equity. *American Journal of Law and Equality*, 1, 167-193. [https://doi.org/10.1162/ajle\\_a\\_00019](https://doi.org/10.1162/ajle_a_00019)
- Muhammad, K. G. (2011). *The Condemnation of Blackness: Race, Crime, and the Making of Modern Urban America*. Harvard University Press. (2011)
- Murakawa, N. (2014). *The First Civil Right: How Liberals Built Prison America*. Oxford University Press.
- O'Neil, C. (2016). ProPublica report: recidivism risk models are racist. *mathbabe*. <https://mathbabe.org/2016/05/24/propublica-report-recidivism-risk-models-are-racist/>
- Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464), 447-453. <https://doi.org/10.1126/science.aax2342>
- Ochigame, R. (2020). The Long History of Algorithmic Fairness. *Phenomenal World*. <https://phenomenalworld.org/analysis/long-history-algorithmic-fairness>
- Ochigame, R., Barabas, C., Dinakar, K., Virza, M., & Ito, J. (2018). Beyond Legitimation: Rethinking Fairness, Interpretability, and Accuracy in Machine Learning. Machine Learning: The Debates workshop at the 35th International Conference on Machine Learning,
- Passi, S., & Barocas, S. (2019). Problem Formulation and Fairness. Proceedings of the Conference on Fairness, Accountability, and Transparency,

- Porrino, C. S. (2017). Attorney General Law Enforcement Directive 2016-6 v3.0: Modification of Directive Establishing Interim Policies, Practices, and Procedures to Implement Criminal Justice Reform Pursuant to P.L. 2015, c. 31. [https://www.nj.gov/lps/dcj/agguide/directives/ag-directive-2016-6\\_v3-0.pdf](https://www.nj.gov/lps/dcj/agguide/directives/ag-directive-2016-6_v3-0.pdf)
- Powles, J., & Nissenbaum, H. (2018). The Seductive Diversion of ‘Solving’ Bias in Artificial Intelligence. *OneZero*. <https://onezero.medium.com/the-seductive-diversion-of-solving-bias-in-artificial-intelligence-890df5e5ef53>
- Raji, I. D., & Buolamwini, J. (2019). Actionable Auditing: Investigating the Impact of Publicly Naming Biased Performance Results of Commercial AI Products. Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, Honolulu, HI, USA.
- Rose, D. R., & Clear, T. R. (1998). Incarceration, Social Capital, and Crime: Implications for Social Disorganization Theory. *Criminology*, 36(3), 441-480.
- Rothstein, R. (2017). *The Color of Law: A Forgotten History of How Our Government Segregated America*. Liveright Publishing Corporation.
- Sampson, R. J., Morenoff, J. D., & Raudenbush, S. (2005). Social Anatomy of Racial and Ethnic Disparities in Violence. *American Journal of Public Health*, 95(2), 224-232. <https://doi.org/10.2105/AJPH.2004.037705>
- Selbst, A. D., boyd, d., Friedler, S. A., Venkatasubramanian, S., & Vertesi, J. (2019). Fairness and Abstraction in Sociotechnical Systems. Proceedings of the Conference on Fairness, Accountability, and Transparency, Atlanta, GA, USA.
- Sunstein, C. R. (2019). Algorithms, Correcting Biases. *Social Research*, 86(2), 499-511.
- Tuttle, C. (2019). Snapping Back: Food Stamp Bans and Criminal Recidivism. *American Economic Journal: Economic Policy*, 11(2), 301-327. <https://doi.org/10.1257/pol.20170490>
- U.S. Supreme Court. (1987). *United States v. Salerno*. 481 U.S. 739.
- Unger, R. M. (2005). *The Left Alternative*. Verso.
- West, S. M. (2020). Redistribution and Rekognition. *Catalyst: Feminism, Theory, Technoscience*, 6(2).
- Wong, P.-H. (2020, 2020/06/01). Democratizing Algorithmic Fairness. *Philosophy & Technology*, 33(2), 225-244. <https://doi.org/10.1007/s13347-019-00355-w>
- Wright, E. O. (2010). *Envisioning Real Utopias*. Verso.
- Ye, T., Johnson, R., Fu, S., Copeny, J., Donnelly, B., Freeman, A., et al. (2019). *Using machine learning to help vulnerable tenants in New York City* Proceedings of the 2nd ACM SIGCAS Conference on Computing and Sustainable Societies, Accra, Ghana. <https://doi.org/10.1145/3314344.3332484>