# Entitlement Theory and Algorithmic Fairness: A Nozickian Framework for Just Decision-Making Systems

## Abstract

This thesis examines the theoretical underpinnings of algorithmic fairness measures through the lens of Nozick's libertarian theory of entitlement justice. While current algorithmic fairness research heavily draws from egalitarian theories of distributive justice, focusing on equality of outcomes across demographic groups, this perspective has neglected libertarian conceptions of fairness that emphasize procedural justice and respect for individual rights. By developing a Nozickian framework for algorithmic fairness, this thesis offers a novel perspective on how algorithmic decision systems might be evaluated against principles of acquisition, transfer, and rectification. The thesis proposes a formal measure of "procedural entitlement fairness" and applies it to case studies in algorithmic lending and hiring decisions. This approach illuminates the philosophical tensions inherent in algorithmic governance and suggests that combining egalitarian and libertarian perspectives may yield more nuanced fairness metrics that respect both group equality and individual procedural rights. The framework developed here demonstrates that Nozick's theories, though rarely applied to algorithmic contexts, offer valuable insights for addressing the normative challenges posed by automated decision-making systems.

## 1. Introduction

The increasing role of algorithms in decision-making processes that affect human lives—from loan approvals and hiring decisions to criminal sentencing and resource allocation—has prompted significant ethical concerns about fairness and justice. As algorithms make or inform consequential decisions, questions about what constitutes "fair" algorithmic behavior have become central to both technical and philosophical discussions (Barocas et al., 2019). The field of algorithmic fairness has emerged as a response to these concerns, with researchers proposing various formal definitions and statistical metrics to evaluate whether an algorithm treats different demographic groups equitably.

These fairness metrics, such as demographic parity, equalized odds, and equal opportunity, have predominantly been conceptualized through an egalitarian lens that focuses on the distribution of outcomes across different groups (Corbett-Davies & Goel, 2018). This approach aligns with dominant theories of distributive justice in the philosophical literature, particularly those that emphasize equality of opportunity or outcomes. However, this focus has led to a notable gap: the relative absence of libertarian perspectives on justice, particularly Robert Nozick's entitlement theory, from discussions of algorithmic fairness.

Nozick's theory of justice, as outlined in *Anarchy, State, and Utopia* (1974),

presents a radically different conception of fairness than most egalitarian theories. Rather than focusing on patterns of distribution (such as equality), Nozick argues that justice is historical and procedural—a distribution is just if it arose through just processes of acquisition and transfer, regardless of the resulting pattern. This perspective has been largely overlooked in algorithmic fairness research, despite its significant influence in political philosophy.

This theoretical gap raises several important questions: How might algorithmic fairness be conceptualized from a Nozickian perspective? What would a measure of algorithmic fairness based on entitlement theory look like? How would such a measure compare to existing egalitarian metrics? And what insights might a libertarian approach offer to the broader discussion of algorithmic ethics?

## 1.1 Research Objectives

This thesis aims to address these questions through the following objectives:

1. To elucidate the philosophical foundations of current algorithmic fairness metrics, highlighting their egalitarian assumptions and limitations.

2. To develop a theoretical framework for algorithmic fairness grounded in Nozick's entitlement theory of justice.

3. To formulate a formal measure or set of measures that can evaluate algorithmic fairness according to Nozickian principles.

4. To apply this framework to case studies in algorithmic decision-making, comparing the results with those derived from traditional egalitarian measures.

5. To explore the philosophical implications and practical challenges of implementing a libertarian approach to algorithmic fairness.

## 1.2 Significance and Contribution

This research contributes to both philosophical and technical discussions of algorithmic fairness in several ways. First, it broadens the philosophical foundations of algorithmic ethics by introducing a perspective that has been underrepresented in the literature. Second, it offers a novel conceptual framework for evaluating algorithmic systems that complements existing approaches. Third, it provides practical measures that can be used by developers and policymakers to assess algorithms from a libertarian perspective. Finally, it illuminates the broader philosophical tensions that underlie debates about fairness in the digital age.

By bringing Nozick's entitlement theory into conversation with algorithmic fairness, this thesis aims to enrich our understanding of what justice means in the context of automated decision-making and to provide new tools for ensuring that algorithms align with diverse conceptions of fairness.

### 1.3 Methodological Approach

This thesis employs a methodology that combines philosophical analysis with formal modeling. The approach includes:

1. Critical analysis of the philosophical foundations of existing algorithmic fairness metrics, with particular attention to their implicit ethical assumptions.

2. Exegetical analysis of Nozick's entitlement theory to extract core principles relevant to algorithmic decision-making.

3. Formal development of metrics that operationalize Nozickian principles in the context of algorithmic systems.

4. Application of these metrics to hypothetical and real-world case studies, with comparative analysis against traditional fairness measures.

5. Philosophical reflection on the implications, limitations, and potential extensions of the proposed framework.

This methodology allows for a rigorous exploration of how libertarian principles might be translated into the technical language of algorithmic fairness while maintaining philosophical integrity.

### 1.4 Thesis Structure

The remainder of this thesis is structured as follows:

Section 2 provides background on algorithmic fairness, including an overview of key metrics and their philosophical underpinnings, as well as an introduction to theories of distributive justice with emphasis on the distinctions between egalitarian and libertarian approaches.

Section 3 develops a Nozickian framework for algorithmic fairness, articulating how the principles of just acquisition, transfer, and rectification can be applied to algorithmic decision-making. This section also proposes formal measures for evaluating algorithmic fairness from a libertarian perspective.

Section 4 applies the developed framework to case studies in algorithmic lending and hiring, comparing the results with those obtained using traditional fairness metrics.

Section 5 discusses the implications of the Nozickian approach, addressing potential objections and limitations, as well as suggesting directions for future research.

Section 6 concludes by summarizing the key contributions of the thesis and reflecting on the broader significance of incorporating diverse philosophical perspectives into algorithmic ethics.

## 2. Background

### 2.1 Algorithmic Fairness: Concepts and Metrics

**2.1.1 The Rise of Algorithmic Decision-Making**  The proliferation of algorithmic decision-making systems across various domains of social life represents one of the most significant technological developments of the 21st century. These systems—ranging from simple rule-based algorithms to complex machine learning models—now influence or determine outcomes in areas such as financial lending (Fuster et al., 2022), hiring (Bogen & Rieke, 2018), criminal justice (Angwin et al., 2016), healthcare resource allocation (Obermeyer et al., 2019), and social service provision (Eubanks, 2018). The increasing reliance on these systems stems from their perceived benefits: efficiency, consistency, scale, and the potential to overcome human biases and limitations.

However, the adoption of algorithmic decision-making has been accompanied by growing concerns about fairness and discrimination. Research has demonstrated that algorithms can reproduce or even amplify existing social biases (Barocas & Selbst, 2016; Noble, 2018). These biases may emerge through various mechanisms, including biased training data, problematic feature selection, proxy discrimination, and the optimization of objectives that disadvantage certain groups (Hellman, 2020). The recognition of these issues has given rise to the field of algorithmic fairness, which seeks to develop techniques for detecting, measuring, and mitigating unfairness in algorithmic systems.

**2.1.2 Formal Definitions of Fairness**  The field of algorithmic fairness has produced numerous formal definitions and statistical metrics to evaluate whether an algorithm treats different demographic groups equitably. These definitions generally fall into several categories, each reflecting different intuitions about what constitutes fair treatment. The most prominent categories include:

**Independence-based fairness (Demographic Parity)**: This definition requires that the algorithm's decision be independent of protected attributes such as race or gender. Formally, if R represents the algorithm's decision (e.g., approve/deny) and A represents a protected attribute (e.g., race), demographic parity requires that $P(R = r \mid A = a) = P(R = r \mid A = a')$ for all values r, a, and a' (Dwork et al., 2012). In other words, the probability of each outcome should be equal across all demographic groups.

**Separation-based fairness (Equalized Odds)**: This definition requires that the algorithm's decision be independent of protected attributes, conditional on the actual outcome Y. Formally, $P(R = r \mid A = a, Y = y) = P(R = r \mid A = a', Y = y)$ for all r, a, a', and y (Hardt et al., 2016). This means that both true positive rates and false positive rates should be equal across demographic groups.

**Equal Opportunity**: A relaxation of equalized odds, this definition only requires equality of true positive rates across groups: $P(R = 1 \mid A = a, Y = 1) =$

P(R = 1 | A = a', Y = 1) for all a and a' (Hardt et al., 2016).

**Sufficiency-based fairness (Calibration)**: This definition requires that the algorithm's predicted probability of an outcome match the actual probability of that outcome within each group. Formally, $P(Y = 1 | R = r, A = a) = P(Y = 1 | R = r, A = a') = r$ for all a and a' (Kleinberg et al., 2017).

**Individual Fairness**: Moving beyond group-level measures, individual fairness requires that similar individuals receive similar treatment. Formally, if $d(x, x')$ represents the distance between individuals x and x' in a relevant metric space, and $d'(R(x), R(x'))$ represents the distance between the algorithm's decisions for these individuals, then $d'(R(x), R(x')) \leq d(x, x')$ (Dwork et al., 2012).

These definitions represent different interpretations of fairness, often stemming from different normative intuitions about what constitutes just treatment. Importantly, research has demonstrated that many of these definitions are mathematically incompatible—it is generally impossible to satisfy all of them simultaneously, except in highly constrained or trivial scenarios (Kleinberg et al., 2017; Chouldechova, 2017). This incompatibility necessitates choices about which fairness criteria to prioritize in any given context, choices that ultimately reflect value judgments about the nature of fairness itself.

**2.1.3 Limitations and Critiques of Current Approaches**  While formal fairness metrics have provided valuable tools for identifying and addressing certain forms of algorithmic bias, they have been subject to various critiques. These critiques can be broadly categorized as follows:

**Technical Limitations**: As noted above, the impossibility results of Kleinberg et al. (2017) and Chouldechova (2017) demonstrate that different formal fairness criteria cannot generally be satisfied simultaneously. This means that choosing to satisfy one criterion necessarily involves trade-offs with others. Additionally, many fairness metrics are sensitive to how the problem is framed and which variables are included in the model (Corbett-Davies & Goel, 2018).

**Contextual Inadequacy**: Fairness metrics often fail to account for the specific social, historical, and institutional contexts in which algorithms operate. What constitutes fair treatment may vary across different domains and social settings, making generic statistical definitions insufficient (Green, 2018; Selbst et al., 2019).

**Procedural Neglect**: Most fairness metrics focus on the distribution of outcomes rather than the processes by which those outcomes are determined. This neglects procedural aspects of fairness, such as transparency, explainability, and the ability to contest decisions (Grgić-Hlača et al., 2018; Binns, 2018).

**Individualist vs. Group-Based Tensions**: There is an inherent tension between group-based fairness metrics, which aim to ensure equality across demographic groups, and individualist approaches, which emphasize treating similar

individuals similarly regardless of group membership (Binns, 2020; Dwork et al., 2012).

**Normative Ambiguity**: The choice of which fairness metric to use involves implicit normative assumptions about what constitutes fair treatment, yet these assumptions are often left unexamined. Different metrics may align with different theories of justice, but this connection is rarely made explicit (Fazelpour & Lipton, 2020; Binns, 2018).

These limitations highlight the need for approaches to algorithmic fairness that are philosophically grounded, contextually sensitive, and capable of addressing both procedural and distributive aspects of justice. This thesis aims to contribute to this effort by exploring how Nozick's entitlement theory might offer a distinctive perspective on algorithmic fairness—one that foregrounds procedural justice and historical processes rather than distributive patterns.

### 2.2 Theories of Distributive Justice

**2.2.1 Egalitarian Theories**  Egalitarian theories of distributive justice are primarily concerned with the way goods, resources, opportunities, or welfare are distributed across society. These theories generally hold that justice requires some form of equality, though they differ in what exactly should be equalized. The main variants of egalitarianism include:

**Strict Egalitarianism**: This view holds that justice requires equality of outcomes—all individuals should receive equal shares of goods or resources. This position is rarely defended in its pure form due to various practical and theoretical objections (such as differences in need, effort, or desert), but it serves as a baseline against which other theories are often contrasted (Nielsen, 1979).

**Luck Egalitarianism**: This approach, associated with philosophers such as Ronald Dworkin (2000), G.A. Cohen (1989), and Richard Arneson (1989), aims to neutralize the effects of brute or unchosen luck while holding individuals responsible for their choices. Justice requires compensating individuals for disadvantages that result from circumstances beyond their control (brute luck) but allows inequalities that result from voluntary choices (option luck).

**Prioritarianism**: Rather than strict equality, prioritarianism holds that justice requires giving priority to improving the condition of those who are worst off. This view, associated with philosophers like Derek Parfit (1997), differs from strict egalitarianism in that it is concerned with absolute rather than relative levels of well-being.

**Rawlsian Justice as Fairness**: John Rawls's (1971) influential theory combines egalitarian and prioritarian elements. His difference principle permits inequalities only if they benefit the least advantaged members of society. Additionally, his principle of fair equality of opportunity requires that positions and

offices be open to all under conditions where individuals with similar talents and motivation have similar prospects for success.

These egalitarian approaches share a focus on the pattern or structure of distribution—they evaluate justice based on how goods or opportunities are distributed at a given time, rather than the processes by which that distribution came about. This pattern-based approach has been influential in shaping algorithmic fairness metrics, particularly those focused on ensuring equality of outcomes or opportunities across demographic groups.

**2.2.2 Libertarian Theories and Nozick's Entitlement Theory** In contrast to egalitarian theories, libertarian approaches to distributive justice focus on individual rights and the processes by which distributions arise, rather than the resulting patterns. The most prominent libertarian theory of justice is Robert Nozick's entitlement theory, articulated in his seminal work *Anarchy, State, and Utopia* (1974).

Nozick's theory is fundamentally historical and procedural—it holds that a distribution is just if it arose through just processes, regardless of the resulting pattern. The theory consists of three main principles:

**The Principle of Justice in Acquisition**: This principle concerns the original acquisition of holdings from an unowned state of nature. Nozick draws on Locke's theory of property acquisition, which holds that one can acquire previously unowned resources by mixing one's labor with them, subject to the "Lockean proviso" that "enough and as good" is left for others (Nozick, 1974, p. 175).

**The Principle of Justice in Transfer**: Once resources have been justly acquired, they may be transferred to others. A transfer is just if it occurs through voluntary exchange, gift, or bequest. Coerced transfers, such as theft or fraud, violate this principle and result in unjust holdings (Nozick, 1974, p. 150-153).

**The Principle of Rectification of Injustice**: This principle addresses how to rectify past injustices in acquisition or transfer. If current holdings resulted from unjust processes, rectification is required to restore a just distribution. Nozick acknowledges the complexity of determining what rectification would entail in real-world scenarios with histories of injustice (Nozick, 1974, p. 152-153).

Nozick's theory explicitly rejects what he calls "patterned" theories of justice—theories that specify that a distribution is to vary along with some natural dimension, such as moral merit, need, or usefulness to society. He argues that maintaining any pattern would require continuous interference with people's ability to freely exchange what they own, thus violating their rights (Nozick, 1974, p. 160-164). His famous "Wilt Chamberlain" example illustrates how liberty upsets patterns: if we start with a just distribution (according to some pattern) and then allow people to freely exchange their resources (e.g., by pay-

ing to watch Wilt Chamberlain play basketball), the resulting distribution will no longer conform to the original pattern. Nozick argues that preventing this outcome would require "continuous interference with people's lives" (Nozick, 1974, p. 163).

Nozick's approach represents a radical departure from egalitarian theories in that it does not evaluate justice based on the resulting distribution of goods or opportunities, but rather on whether that distribution arose through processes that respect individual rights to property and free exchange. This procedural focus offers a distinctive perspective that has been largely absent from discussions of algorithmic fairness.

**2.2.3 Critiques and Counterarguments**  Both egalitarian and libertarian theories have been subject to extensive critique. Key criticisms of egalitarian theories include concerns about the restriction of liberty, problems with determining the appropriate "currency" of equality (resources, welfare, capabilities, etc.), and questions about the moral relevance of equality itself as opposed to sufficiency or priority (Frankfurt, 1987; Parfit, 1997).

Critiques of Nozick's entitlement theory have been equally robust. Major objections include:

**The Problem of Initial Acquisition**: Critics argue that Nozick's account of just acquisition is underdeveloped and that the Lockean proviso may be impossible to satisfy in a world of finite resources (Cohen, 1995; Otsuka, 2003).

**Historical Injustice**: Given the extensive history of injustice in acquisition and transfer (including colonialism, slavery, and theft), critics argue that almost all current holdings are tainted by past injustice, rendering Nozick's theory practically inapplicable without massive redistribution (Kymlicka, 2002; Waldron, 1992).

**Neglect of Need**: Critics contend that Nozick's theory fails to account for basic human needs and allows for distributions that leave some individuals in desperate poverty while others have vast wealth (Cohen, 1995).

**Self-Ownership Without Substance**: Critics like G.A. Cohen (1995) argue that Nozick's emphasis on self-ownership is hollow if individuals lack access to external resources needed to make meaningful use of their talents and abilities.

Despite these critiques, Nozick's theory remains influential and offers valuable insights into procedural aspects of justice that are often neglected in egalitarian approaches. The tension between procedural and distributive conceptions of justice—between process and pattern—is central to many debates in political philosophy and, as this thesis argues, equally relevant to discussions of algorithmic fairness.

### 2.3 The Intersection of Algorithmic Fairness and Distributive Justice

**2.3.1 Current Perspectives**  Recent scholarship has begun to explore the connections between algorithmic fairness and theories of distributive justice, primarily focusing on how different fairness metrics align with different philosophical conceptions of fairness (Binns, 2018; Fazelpour & Lipton, 2020; Lee & Floridi, 2021).

This work has revealed that many common fairness metrics implicitly adopt egalitarian perspectives. For example:

- Demographic parity aligns with a strict egalitarian view that different groups should receive equal outcomes regardless of other factors.
- Equal opportunity metrics reflect Rawlsian concerns with fair equality of opportunity.
- Calibration and predictive parity metrics can be linked to meritocratic conceptions of desert-based justice.

These connections have helped clarify the normative assumptions underlying different technical approaches to fairness and have highlighted the need for contextual judgment in selecting appropriate metrics for specific applications (Lee & Floridi, 2021).

However, as noted in the introduction, libertarian perspectives—particularly Nozick's entitlement theory—have been largely absent from these discussions. This absence is notable given the significant influence of libertarian thinking in debates about technology, markets, and regulation more broadly (Thierer, 2016).

**2.3.2 The Gap: Libertarian Perspectives on Algorithmic Fairness**
The omission of libertarian perspectives from discussions of algorithmic fairness represents a significant gap in the literature. This gap may stem from several factors:

1. The apparent incompatibility between libertarian emphasis on individual rights and process versus the focus on group-level statistics in most fairness metrics.

2. The challenge of translating historical and procedural conceptions of justice into formal, quantifiable measures applicable to algorithms.

3. The dominant egalitarian orientation of much research on discrimination and bias, reflecting broader trends in social justice scholarship.

4. The difficulty of applying Nozick's principles to algorithmic contexts where the concepts of "acquisition" and "transfer" do not have obvious analogues.

Despite these challenges, a Nozickian perspective on algorithmic fairness offers potential insights that complement existing approaches. By focusing on the processes by which algorithmic decisions are made, rather than solely on the

resulting distributions, such a perspective might address some of the procedural concerns that have been raised about current fairness metrics (Grgić-Hlača et al., 2018; Binns, 2018). Additionally, a libertarian approach might help navigate tensions between group fairness and individual treatment that have proven challenging for egalitarian metrics (Binns, 2020).

In the following section, I develop a framework for algorithmic fairness grounded in Nozick's entitlement theory, addressing the challenges of translating his historical and procedural conception of justice into the context of algorithmic decision-making.

## 3. A Nozickian Framework for Algorithmic Fairness

### 3.1 Conceptual Foundations

**3.1.1 Core Elements of Nozick's Theory Relevant to Algorithms** To develop a Nozickian framework for algorithmic fairness, we must first identify which elements of Nozick's entitlement theory are most relevant to algorithmic decision-making contexts. While the direct translation of concepts like "acquisition" and "transfer" from physical property to algorithmic decisions is not straightforward, several core aspects of Nozick's theory provide valuable starting points:

**Procedural Justice**: For Nozick, justice is fundamentally about processes rather than patterns. A just distribution is one that arises through just processes, regardless of the resulting pattern. This focus on procedural justice can be applied to algorithmic contexts by emphasizing the procedures through which algorithms make decisions, rather than solely the distribution of outcomes across groups.

**Historical Dimension**: Nozick's theory is historical—the justice of current holdings depends on the history of how they came to be. In algorithmic contexts, this suggests attention to the historical processes that generate the data used by algorithms and the historical context in which algorithmic systems operate.

**Rights and Consent**: Central to Nozick's theory is respect for individual rights, particularly property rights and the right to engage in voluntary exchanges. In algorithmic contexts, this might translate to concerns about consent, data ownership, and individuals' rights regarding how their information is used in decision-making processes.

**Non-Interference**: Nozick argues against "continuous interference" to maintain patterns of distribution. In algorithmic contexts, this might suggest skepticism toward extensive interventions in algorithm design solely to achieve particular distributive outcomes, especially if these interventions might violate individual rights or distort market processes.

**Rectification**: Nozick acknowledges the need to rectify past injustices in acquisition and transfer. In algorithmic contexts, this principle might apply to

addressing historical biases in data or correcting for past discriminatory practices that may be perpetuated by algorithms.

These elements provide a foundation for thinking about algorithmic fairness from a Nozickian perspective. However, translating them into concrete principles for algorithmic design and evaluation requires addressing several conceptual challenges.

**3.1.2 Translating Nozickian Concepts to Algorithmic Contexts**  The translation of Nozick's principles to algorithmic contexts requires addressing several key questions:

**What constitutes "just acquisition" in algorithmic decision-making?** In Nozick's theory, just acquisition concerns how unowned resources come to be legitimately owned. In algorithmic contexts, this might relate to how data is collected and how algorithmic models are developed. Just acquisition could involve obtaining data with proper consent, ensuring transparency about how data will be used, and developing algorithms in ways that respect the rights and autonomy of affected individuals.

**What constitutes "just transfer" in algorithmic decisions?** For Nozick, just transfer involves voluntary exchanges free from coercion or fraud. In algorithmic contexts, this might relate to how decisions are made and communicated. Just transfer could involve ensuring that algorithmic decisions are transparent, that individuals understand the basis for decisions affecting them, and that they have meaningful opportunities to contest or appeal decisions they believe are unjust.

**How does the principle of rectification apply to algorithms?** Nozick's principle of rectification addresses how to correct past injustices in acquisition or transfer. In algorithmic contexts, this might involve correcting for historical biases in data, addressing past discriminatory practices that may be perpetuated by algorithms, or providing remedies for individuals who have been harmed by unjust algorithmic decisions.

**What rights do individuals have in relation to algorithmic systems?** Nozick's theory is fundamentally concerned with individual rights, particularly property rights. In algorithmic contexts, relevant rights might include data ownership rights, rights to consent to how one's data is used, rights to explanation or contestation of algorithmic decisions, and rights to compensation for harms caused by algorithmic systems.

**How do we balance procedural justice with concerns about discriminatory outcomes?** While Nozick's theory focuses on procedural justice rather than distributive patterns, there may be cases where seemingly just procedures lead to outcomes that appear discriminatory. Addressing this tension requires careful consideration of how procedural and distributive concerns can be balanced within a broadly libertarian framework.

Addressing these questions will help us develop concrete principles for evaluating algorithmic fairness from a Nozickian perspective. In the following sections, I propose such principles and develop formal metrics based on them.

### 3.2 Principles of Nozickian Algorithmic Fairness

Building on the conceptual foundations outlined above, I propose the following principles for evaluating algorithmic fairness from a Nozickian perspective:

**3.2.1 Principle of Just Data Acquisition**  An algorithm satisfies the principle of just data acquisition if:

1. All data used to train and validate the algorithm was collected with the informed consent of the individuals to whom the data pertains, or from legitimately public sources.
2. Individuals retained meaningful control over their data, including rights to access, correct, delete, or restrict the use of their personal information.
3. The collection and use of data respected relevant privacy rights and did not involve deception or coercion.
4. The Lockean proviso is satisfied: the data collection did not deprive others of essential information resources or create significant informational asymmetries that undermine individual autonomy.

This principle translates Nozick's concern with just initial acquisition to the context of data collection and algorithm development. It emphasizes consent, control, and respect for individual rights over personal information.

**3.2.2 Principle of Just Algorithmic Processing**  An algorithm satisfies the principle of just algorithmic processing if:

1. The algorithm makes decisions based on factors that individuals have meaningfully influenced through their voluntary choices and actions, rather than immutable characteristics or circumstances beyond their control.
2. The algorithm's decision-making process is transparent and intelligible to affected individuals, allowing them to understand how their actions and choices influence decisions about them.
3. The algorithm does not violate individuals' rights or entitlements established through just processes (e.g., legal contracts, earned credentials, legitimate expectations based on past performance).
4. The algorithm's operation does not involve fraud, deception, or other forms of involuntary transfer (e.g., making decisions based on criteria different from those publicly claimed).

This principle applies Nozick's focus on voluntary transfers and respect for established entitlements to the context of algorithmic decision-making. It emphasizes transparency, agency, and respect for rightfully acquired claims.

### 3.2.3 Principle of Algorithmic Rectification

An algorithm satisfies the principle of algorithmic rectification if:

1. It includes mechanisms to identify and correct for past injustices that may be perpetuated through algorithmic decisions, particularly those involving violations of just acquisition or transfer.
2. It provides meaningful opportunities for individuals to contest decisions they believe are unjust and to seek appropriate remedies.
3. It includes processes for updating and improving the algorithm in response to identified instances of injustice or rights violations.
4. When historical data reflects past discriminatory practices or unjust social arrangements, the algorithm incorporates appropriate adjustments to prevent the perpetuation of these injustices.

This principle applies Nozick's principle of rectification to algorithmic contexts, acknowledging that historical injustices may require correction to achieve a just state of affairs. It recognizes that while procedural justice is paramount, there may be cases where past procedural violations necessitate interventions to restore justice.

### 3.2.4 Meta-Principle: Minimal Interference with Voluntary Exchanges

In addition to the three substantive principles outlined above, a Nozickian approach to algorithmic fairness would include a meta-principle regarding the role of regulation and intervention:

Interventions in algorithmic design and operation should be limited to those necessary to ensure just acquisition, processing, and rectification, and should minimize interference with voluntary exchanges and individual rights. This principle reflects Nozick's concern with "continuous interference" to maintain particular distributive patterns.

This meta-principle does not preclude all regulation or intervention—Nozick himself acknowledges the legitimacy of state action to protect rights and enforce just acquisition and transfer. However, it suggests skepticism toward extensive interventions solely to achieve particular distributional outcomes, especially if these interventions might infringe on individual rights or distort market processes.

### 3.3 Formalizing Nozickian Fairness Metrics

Translating the principles outlined above into formal metrics that can be applied to algorithmic systems presents a significant challenge. Unlike traditional fairness metrics, which typically focus on statistical properties of an algorithm's outputs, Nozickian metrics must capture procedural aspects of justice that are not easily quantified. Nevertheless, I propose the following formal approaches to measuring algorithmic fairness from a Nozickian perspective:

**3.3.1 Procedural Entitlement Fairness (PEF)**   I propose a metric called Procedural Entitlement Fairness (PEF) that evaluates the extent to which an algorithm's decisions respect individuals' rightfully acquired entitlements. Let's define:

- $E_i$ = the set of rightfully acquired entitlements of individual $i$
- $D_i$ = the decision made by the algorithm regarding individual $i$
- $R(E_i, D_i)$ = a function that measures the extent to which decision $D_i$ respects the entitlements $E_i$

Then, Procedural Entitlement Fairness (PEF) for a population of $n$ individuals can be defined as:

$$PEF = \frac{1}{n} \sum_{i=1}^{n} R(E_i, D_i)$$

Where $R(E_i, D_i)$ is normalized to a value between 0 and 1, with 1 indicating full respect for entitlements and 0 indicating complete violation.

The challenge in implementing this metric lies in defining $E_i$ and $R(E_i, D_i)$ for specific contexts. In a lending context, for example, $E_i$ might include an individual's credit history, repayment record, and legitimate expectations based on past financial behavior, while $R(E_i, D_i)$ might measure how well the lending decision aligns with these established entitlements.

**3.3.2 Consent and Transparency Index (CTI)**   The Consent and Transparency Index (CTI) measures the extent to which an algorithmic system respects principles of just acquisition and transparent processing. Let's define:

- $C_i$ = a binary indicator of whether individual $i$'s data was collected with proper consent (1) or not (0)
- $T_i$ = a measure of the transparency of the algorithm's decision process for individual $i$, normalized to a value between 0 and 1
- $w_c$ and $w_t$ = weights assigned to the consent and transparency components, respectively, with $w_c + w_t = 1$

Then, the Consent and Transparency Index (CTI) can be defined as:

$$CTI = \frac{1}{n} \sum_{i=1}^{n} (w_c \cdot C_i + w_t \cdot T_i)$$

This metric captures key aspects of just acquisition (consent) and just processing (transparency), reflecting the procedural focus of Nozick's theory.

### 3.3.3 Rectification Responsiveness Measure (RRM)

The Rectification Responsiveness Measure (RRM) evaluates an algorithm's capacity to identify and correct for past injustices that may be perpetuated through its decisions. Let's define:

- $I$ = the set of identified instances where the algorithm's decisions may perpetuate past injustices
- $C_j$ = a measure of the adequacy of the correction implemented for instance $j \in I$, normalized to a value between 0 and 1
- $A$ = a measure of the algorithm's accessibility to contestation and appeal, normalized to a value between 0 and 1
- $w_c$ and $w_a$ = weights assigned to the correction and accessibility components, respectively, with $w_c + w_a = 1$

Then, the Rectification Responsiveness Measure (RRM) can be defined as:

$RRM = w_c \cdot \left( \frac{1}{|I|} \sum_{j \in I} C_j \right) + w_a \cdot A$

This metric captures the algorithm's ability to identify and correct for past injustices, as well as its accessibility to contestation by affected individuals, reflecting Nozick's principle of rectification.

### 3.3.4 Choice Sensitivity Ratio (CSR)

The Choice Sensitivity Ratio (CSR) measures the extent to which an algorithm's decisions are sensitive to factors that individuals can influence through their voluntary choices, as opposed to immutable characteristics or circumstances beyond their control. Let's define:

- $F_v$ = the set of features used by the algorithm that reflect voluntary choices
- $F_i$ = the set of features used by the algorithm that reflect immutable characteristics or circumstances beyond individual control
- $I_v$ = the importance of voluntary features in the algorithm's decisions, measured by their collective influence on the algorithm's output
- $I_i$ = the importance of immutable features in the algorithm's decisions, measured similarly

Then, the Choice Sensitivity Ratio (CSR) can be defined as:

$CSR = \frac{I_v}{I_v + I_i}$

This ratio ranges from 0 to 1, with higher values indicating greater sensitivity to voluntary choices, which aligns with Nozick's emphasis on individual agency and responsibility.

### 3.4 Integrated Nozickian Fairness Framework

While the individual metrics proposed above capture specific aspects of Nozickian fairness, an integrated framework is needed to evaluate algorithmic systems comprehensively. I propose an integrated Nozickian Fairness Score (NFS) that combines the four metrics:

$$NFS = w_{PEF} \cdot PEF + w_{CTI} \cdot CTI + w_{RRM} \cdot RRM + w_{CSR} \cdot CSR$$

Where $w_{PEF}$, $w_{CTI}$, $w_{RRM}$, and $w_{CSR}$ are weights assigned to each component, with $w_{PEF} + w_{CTI} + w_{RRM} + w_{CSR} = 1$.

The weights can be adjusted based on the specific context and the relative importance of different aspects of Nozickian fairness in that context. For example, in contexts where historical injustices are particularly salient, greater weight might be given to RRM, while in contexts where individual agency is paramount, CSR might receive greater weight.

This integrated framework allows for a comprehensive evaluation of algorithmic systems from a Nozickian perspective, focusing on procedural justice rather than distributive patterns. Importantly, it does not preclude the possibility of unequal outcomes across demographic groups—if these inequalities arise through just processes, they are not considered unjust from a Nozickian perspective.

**3.4.1 Comparison with Traditional Fairness Metrics** The Nozickian fairness framework differs significantly from traditional fairness metrics in several key respects:

**Focus on Process vs. Outcome**: Traditional metrics like demographic parity and equalized odds focus on the distribution of outcomes across demographic groups. In contrast, the Nozickian framework focuses on the processes by which decisions are made, emphasizing consent, transparency, respect for entitlements, and sensitivity to voluntary choices.

**Individual vs. Group Level**: Traditional metrics typically operate at the group level, comparing outcomes across demographic categories. The Nozickian framework primarily operates at the individual level, evaluating how algorithmic decisions respect each individual's rights and entitlements.

**Historical vs. Snapshot View**: Traditional metrics typically take a snapshot view, evaluating the current distribution of outcomes without reference to how that distribution came about. The Nozickian framework takes a historical view, considering how data was acquired and how past injustices might require rectification.

**Acceptance of Unequal Outcomes**: Traditional metrics often aim for some form of equality or parity in outcomes across groups. The Nozickian framework accepts potentially unequal outcomes if they arise through just processes of acquisition, transfer, and rectification.

These differences illustrate the distinct perspective that a Nozickian approach brings to discussions of algorithmic fairness. Rather than asking whether an algorithm produces equal outcomes across groups, it asks whether the algorithm respects individual rights, operates transparently, and provides appropriate mechanisms for rectification when needed.

## 4. Application to Case Studies

To illustrate how the Nozickian framework for algorithmic fairness might be applied in practice, I examine two hypothetical case studies: algorithmic lending decisions and automated hiring systems. These applications demonstrate the framework's practical utility while highlighting how it differs from traditional fairness approaches.

### 4.1 Case Study 1: Algorithmic Lending Decisions

**4.1.1 Context** Consider a financial institution that uses an algorithmic system to evaluate loan applications. The system uses various data points—including credit history, income, employment stability, and debt-to-income ratio—to predict an applicant's likelihood of repaying a loan. Based on this prediction, the algorithm either approves or denies the loan application, or offers different interest rates to different applicants.

Traditional fairness metrics might focus on whether the algorithm's decisions result in equal approval rates or similar interest rates across demographic groups such as race or gender. From an egalitarian perspective, substantial disparities in outcomes might be considered evidence of unfairness, even if the algorithm accurately predicts repayment likelihood.

Let us examine this scenario through the lens of Nozickian fairness, applying the metrics developed in Section 3.

**4.1.2 Just Data Acquisition** To evaluate the principle of just data acquisition, we would consider:

- Was the data used to train and validate the algorithm collected with informed consent?
- Do individuals retain control over their financial data, with rights to access, correct, and restrict its use?
- Does the data collection respect privacy rights and avoid deception or coercion?
- Does the data collection create significant informational asymmetries that undermine individual autonomy?

If the algorithm uses credit bureau data, we would examine whether individuals provided informed consent for this use and whether they have meaningful control over their credit information. If the algorithm scrapes social media or other personal data without explicit consent, this would violate the principle of just acquisition.

Applying the Consent and Transparency Index (CTI), we might find that while traditional credit data is collected with some form of consent (though perhaps not fully informed consent), the use of alternative data sources might lack proper consent, resulting in a lower CTI score.

**4.1.3 Just Algorithmic Processing**   To evaluate the principle of just algorithmic processing, we would consider:

- Does the algorithm base decisions primarily on factors that individuals have meaningfully influenced through their voluntary choices (e.g., payment history, debt management) rather than immutable characteristics or circumstances beyond their control (e.g., race, family background)?
- Is the algorithm's decision-making process transparent and intelligible to loan applicants?
- Does the algorithm respect legitimate expectations based on an individual's credit history and financial behavior?
- Does the algorithm operate as described, without hidden factors or deceptive practices?

Applying the Procedural Entitlement Fairness (PEF) metric, we would assess whether individuals with similar credit histories and financial behaviors receive similar treatment, regardless of demographic characteristics. The Choice Sensitivity Ratio (CSR) would measure the extent to which loan decisions are based on voluntary financial behaviors rather than immutable characteristics or circumstances beyond individual control.

**4.1.4 Algorithmic Rectification**   To evaluate the principle of algorithmic rectification, we would consider:

- Does the algorithm account for historical discrimination in lending that may be reflected in credit histories?
- Are there meaningful opportunities for individuals to contest loan denials or unfavorable terms?
- Does the algorithm improve over time based on identified instances of unjust decisions?
- Are there mechanisms to address cases where historical data reflects past discriminatory practices?

Applying the Rectification Responsiveness Measure (RRM), we would assess whether the lending algorithm includes mechanisms to identify and correct for past discriminatory lending practices that might be perpetuated through algorithmic decisions.

**4.1.5 Integrated Evaluation**   Combining these assessments into the integrated Nozickian Fairness Score (NFS), we might find that a lending algorithm scores highly on some dimensions (e.g., basing decisions on voluntary financial behaviors) but poorly on others (e.g., limited transparency or weak rectification mechanisms).

Importantly, from a Nozickian perspective, disparities in loan approval rates across demographic groups would not necessarily indicate unfairness if these disparities arose through just processes. If individuals from different groups have different credit histories due to their voluntary financial choices, and the

algorithm bases decisions on these histories in a transparent and consistent manner, the resulting disparities would not be considered unjust.

However, if disparities arise from historical injustices in acquisition or transfer (e.g., past discriminatory lending practices that affected credit histories), the principle of rectification would require appropriate adjustments to prevent the perpetuation of these injustices. This might involve specific interventions to account for the effects of past discrimination while still respecting individual financial behavior.

**4.1.6 Comparison with Egalitarian Approaches** The Nozickian approach to evaluating the lending algorithm differs markedly from traditional egalitarian approaches. While an egalitarian approach might focus on achieving similar approval rates or loan terms across demographic groups, the Nozickian approach focuses on the processes by which lending decisions are made.

This difference is particularly evident in cases where disparities in outcomes result from differences in credit histories that reflect voluntary financial choices. An egalitarian approach might view such disparities as problematic and advocate for interventions to achieve more equal outcomes. A Nozickian approach would accept these disparities as just if they arose through processes that respected principles of just acquisition, transfer, and rectification.

However, in cases where disparities reflect historical injustices rather than voluntary choices, both approaches might support interventions, though for different reasons. An egalitarian approach would seek to reduce disparities for their own sake, while a Nozickian approach would seek to rectify past violations of just acquisition or transfer.

This comparison highlights both the distinctions and potential points of convergence between egalitarian and libertarian approaches to algorithmic fairness. While they differ in their fundamental orientations—pattern versus process—they may sometimes support similar interventions, particularly in cases involving historical injustice.

**4.2 Case Study 2: Automated Hiring Systems**

**4.2.1 Context** Consider a company that uses an automated system to screen job applicants. The system analyzes resumes, cover letters, and possibly video interviews to predict which candidates are likely to succeed in the role. Based on these predictions, the system either advances candidates to the next stage of the hiring process or rejects their applications.

Traditional fairness metrics might focus on whether the algorithm's decisions result in similar selection rates across demographic groups, with disparities potentially viewed as evidence of discrimination. Let us examine this scenario through the lens of Nozickian fairness.

**4.2.2 Just Data Acquisition**   To evaluate the principle of just data acquisition, we would consider:

- Was the training data (e.g., past hiring decisions, employee performance reviews) collected with informed consent from the individuals involved?
- Do job applicants consent to having their applications processed by an automated system?
- Is the data collection process transparent about how applicant information will be used?
- Does the collection and use of applicant data respect privacy rights and avoid deception?

Applying the Consent and Transparency Index (CTI), we would assess the extent to which the system operates with proper consent and transparency. If the system uses social media data or other personal information without explicit consent, this would lower the CTI score.

**4.2.3 Just Algorithmic Processing**   To evaluate the principle of just algorithmic processing, we would consider:

- Does the algorithm base decisions primarily on factors that reflect candidates' voluntary choices and actions (e.g., education, skills development, work experience) rather than immutable characteristics or circumstances beyond their control (e.g., race, gender, socioeconomic background)?
- Is the algorithm's decision-making process transparent and intelligible to job applicants?
- Does the algorithm respect legitimate qualifications and credentials that candidates have rightfully acquired?
- Does the algorithm operate as described, without hidden factors or deceptive practices?

Applying the Choice Sensitivity Ratio (CSR), we would measure the extent to which the hiring algorithm's decisions are sensitive to factors that reflect candidates' voluntary choices rather than immutable characteristics. The Procedural Entitlement Fairness (PEF) would assess whether candidates with similar qualifications receive similar treatment, regardless of demographic characteristics.

**4.2.4 Algorithmic Rectification**   To evaluate the principle of algorithmic rectification, we would consider:

- Does the algorithm account for historical discrimination in hiring that may be reflected in the training data?
- Are there meaningful opportunities for candidates to contest rejections they believe are unjust?
- Does the system improve over time based on identified instances of unfair decisions?
- Are there mechanisms to address cases where historical hiring data reflects past discriminatory practices?

Applying the Rectification Responsiveness Measure (RRM), we would assess whether the hiring algorithm includes mechanisms to identify and correct for past discriminatory hiring practices that might be perpetuated through algorithmic decisions.

**4.2.5 Integrated Evaluation**  Combining these assessments into the integrated Nozickian Fairness Score (NFS), we might find that a hiring algorithm scores highly on some dimensions (e.g., respecting candidates' rightfully acquired credentials) but poorly on others (e.g., limited mechanisms for contestation or rectification).

From a Nozickian perspective, disparities in selection rates across demographic groups would not necessarily indicate unfairness if these disparities arose through just processes. If individuals from different groups have different qualifications due to their voluntary choices, and the algorithm bases decisions on these qualifications in a transparent and consistent manner, the resulting disparities would not be considered unjust.

However, if disparities arise from historical injustices that affected educational or career opportunities, the principle of rectification would require appropriate adjustments. This might involve specific interventions to account for the effects of past discrimination while still respecting individual qualifications and achievements.

**4.2.6 Comparison with Egalitarian Approaches**  The Nozickian approach to evaluating the hiring algorithm differs from traditional egalitarian approaches in its focus on process rather than outcome patterns. While an egalitarian approach might focus on achieving similar selection rates across demographic groups, the Nozickian approach focuses on the processes by which hiring decisions are made.

This difference is particularly evident in cases where disparities in outcomes result from differences in qualifications that reflect voluntary educational and career choices. An egalitarian approach might view such disparities as problematic and advocate for interventions to achieve more equal outcomes. A Nozickian approach would accept these disparities as just if they arose through processes that respected principles of just acquisition, transfer, and rectification.

However, in cases where disparities reflect historical injustices rather than voluntary choices, both approaches might support interventions, though for different reasons. An egalitarian approach would seek to reduce disparities for their own sake, while a Nozickian approach would seek to rectify past violations of just acquisition or transfer.

This comparison highlights how different philosophical perspectives lead to different evaluations of algorithmic fairness, with potentially significant implications for how automated systems are designed, deployed, and regulated.

# 5. Discussion and Implications

## 5.1 Philosophical Implications

The Nozickian framework for algorithmic fairness developed in this thesis has several important philosophical implications for how we conceptualize and evaluate fairness in automated decision-making systems.

### 5.1.1 Pluralism in Conceptions of Algorithmic Fairness

The development of a libertarian approach to algorithmic fairness alongside existing egalitarian approaches highlights the inherent pluralism in conceptions of fairness. Different philosophical traditions offer distinct lenses through which to evaluate algorithmic systems, with different normative commitments and priorities. This pluralism suggests that there is no single, universally applicable definition of "fairness" that can be encoded into algorithms.

Instead, the choice of fairness metrics and frameworks should be recognized as inherently value-laden, reflecting particular philosophical commitments about what constitutes just treatment. This recognition calls for greater transparency about the normative assumptions underlying different approaches to algorithmic fairness and for more explicit engagement with the philosophical foundations of these approaches.

### 5.1.2 Procedural vs. Distributive Justice in Algorithmic Contexts

The Nozickian framework foregrounds procedural aspects of justice that are often neglected in traditional fairness metrics. While most existing metrics focus on the distribution of outcomes across demographic groups, the Nozickian approach emphasizes the processes by which these outcomes are determined, including issues of consent, transparency, respect for entitlements, and rectification of past injustices.

This procedural focus offers a valuable complement to distributive approaches, highlighting aspects of algorithmic systems that may be morally significant regardless of their distributional effects. It suggests that even algorithms that produce equal outcomes across groups may be procedurally unjust if they violate principles of consent, transparency, or respect for rightfully acquired entitlements.

### 5.1.3 The Role of History in Algorithmic Justice

Nozick's entitlement theory is fundamentally historical, evaluating the justice of current holdings based on how they came to be. Similarly, the Nozickian framework for algorithmic fairness emphasizes the historical dimension of justice, considering how data was acquired, how past injustices might affect current algorithmic decisions, and what rectification might be required.

This historical perspective contrasts with the ahistorical approach of most traditional fairness metrics, which evaluate the current distribution of outcomes

without reference to how that distribution came about. The Nozickian framework suggests that historical context matters for algorithmic fairness—that we cannot evaluate the justice of algorithmic decisions solely by looking at their immediate effects without considering the historical processes that shaped the data and social contexts in which these algorithms operate.

**5.1.4 Individual Rights vs. Group Fairness**  The Nozickian framework prioritizes individual rights and entitlements over group-level patterns of distribution. This individualist orientation contrasts with the group-based focus of most traditional fairness metrics, which compare outcomes across demographic categories. The framework suggests that focusing exclusively on group-level statistics may obscure important aspects of justice at the individual level.

This tension between individual and group perspectives on justice reflects broader debates in political philosophy about the proper subjects of justice. The Nozickian framework suggests that individuals, rather than groups, should be the primary subjects of justice in algorithmic contexts, while acknowledging that group-level patterns may be relevant evidence for identifying potential violations of individual rights.

## 5.2 Practical Implications

Beyond its philosophical significance, the Nozickian framework for algorithmic fairness has several practical implications for how algorithmic systems are designed, deployed, and regulated.

**5.2.1 Implications for Algorithm Design**  The Nozickian framework suggests several principles that should guide the design of algorithmic systems:

**Transparency and Explainability**: Algorithms should be designed to be as transparent and explainable as possible, allowing individuals to understand how decisions affecting them are made and how their actions and choices influence these decisions.

**Consent-Based Data Practices**: The collection and use of data for algorithmic systems should be based on informed consent, with individuals retaining meaningful control over their personal information.

**Respect for Entitlements**: Algorithms should respect individuals' rightfully acquired entitlements, such as credentials, qualifications, or legitimate expectations based on past performance.

**Contestability and Appeals**: Algorithmic systems should include mechanisms for individuals to contest decisions they believe are unjust and to seek appropriate remedies.

**Historical Awareness**: Algorithm designers should be aware of historical injustices that may be reflected in training data and should incorporate appropriate mechanisms to prevent the perpetuation of these injustices.

These design principles differ from those that might be derived from egalitarian approaches to fairness, which might focus more on achieving particular distributional patterns across demographic groups. The Nozickian principles emphasize procedural aspects of algorithmic systems rather than their distributional effects.

**5.2.2 Implications for Regulation and Policy**  The Nozickian framework also has implications for how algorithmic systems are regulated and governed:

**Procedural Requirements**: Regulations might focus on procedural requirements such as transparency, explainability, and contestability rather than mandating particular distributional outcomes.

**Minimal Interference**: Consistent with Nozick's meta-principle of minimal interference, regulations should be limited to those necessary to protect individual rights and ensure just processes, avoiding extensive interventions solely to achieve particular distributional patterns.

**Rights-Based Approach**: Regulations might adopt a rights-based approach, focusing on protecting individuals' rights to consent, explanation, contestation, and rectification rather than imposing specific fairness metrics.

**Context-Sensitivity**: Different contexts may call for different regulatory approaches, depending on the specific rights and entitlements at stake and the historical background against which algorithmic systems operate.

These regulatory implications contrast with approaches that might mandate specific distributive outcomes, such as requiring equal approval rates across demographic groups. The Nozickian approach suggests a more procedural and rights-based orientation to algorithmic governance.

**5.2.3 Complementarity with Egalitarian Approaches**  Despite the differences between Nozickian and egalitarian approaches to algorithmic fairness, they need not be seen as mutually exclusive. In many contexts, procedural and distributive concerns may complement each other, with procedural requirements supporting more equitable distributions and distributional patterns serving as evidence of procedural fairness or unfairness.

For example, in contexts with histories of discrimination, significant disparities in algorithmic outcomes across demographic groups might serve as prima facie evidence of procedural unfairness, prompting investigation into whether principles of just acquisition, transfer, or rectification have been violated. Conversely, ensuring procedural fairness through transparency, consent, and respect for entitlements might naturally lead to more equitable distributions of outcomes.

This complementarity suggests that a comprehensive approach to algorithmic fairness might draw on both Nozickian and egalitarian perspectives, attending to both procedural and distributive aspects of justice. Such a pluralistic approach

would recognize the complex and multifaceted nature of fairness and the need for multiple lenses through which to evaluate algorithmic systems.

### 5.3 Limitations and Challenges

While the Nozickian framework offers valuable insights for algorithmic fairness, it also faces several limitations and challenges that must be acknowledged.

**5.3.1 Practical Implementation Challenges**  Several practical challenges arise when attempting to implement the Nozickian metrics proposed in this thesis:

**Defining Entitlements**: The Procedural Entitlement Fairness (PEF) metric requires defining what constitutes a "rightfully acquired entitlement" in specific contexts, which may be complex and contestable.

**Measuring Voluntary Choice**: The Choice Sensitivity Ratio (CSR) requires distinguishing between factors that reflect voluntary choices and those that reflect immutable characteristics or circumstances beyond individual control, which is often difficult in practice.

**Identifying Past Injustices**:  The Rectification Responsiveness Measure (RRM) requires identifying instances where algorithmic decisions may perpetuate past injustices, which presupposes agreement about what constitutes historical injustice.

**Balancing Components**: The integrated Nozickian Fairness Score (NFS) requires assigning weights to different components, which involves value judgments about their relative importance.

These challenges highlight the need for context-specific implementation of the Nozickian framework, with careful attention to the particular rights, entitlements, and historical injustices relevant to each domain.

**5.3.2 Theoretical Limitations**  Beyond practical challenges, the Nozickian framework also faces several theoretical limitations:

**Initial Acquisition Problem**: As with Nozick's original theory, the framework faces challenges regarding what constitutes just initial acquisition, particularly in the context of data that may have been collected under conditions of unequal power or information asymmetry.

**Historical Complexity**: The historical dimension of the framework requires addressing complex questions about past injustices and appropriate rectification, which may be difficult to resolve definitively.

**Tension with Structural Perspectives**: The framework's focus on individual rights and procedural justice may not adequately address structural forms of injustice that operate without clear violations of individual rights.

**Market Assumptions**: Like Nozick's theory, the framework may rely on assumptions about the fairness of market processes that are contested by critics who point to market failures, power imbalances, and structural inequalities.

These theoretical limitations suggest the need for ongoing philosophical reflection on the foundations of the Nozickian approach and its application to algorithmic contexts.

**5.3.3 Contextual Limitations**   Finally, the Nozickian framework may be more applicable in some contexts than others:

**Market Contexts**: The framework may be most applicable in market contexts where voluntary exchanges and property rights are central, such as lending or hiring.

**Public Services**: The framework may be less applicable in contexts involving public services or goods to which individuals may have claims based on need or citizenship rather than on entitlements acquired through market exchanges.

**Fundamental Rights**: In contexts involving fundamental rights or basic needs, egalitarian or sufficientarian approaches might take precedence over the procedural focus of the Nozickian framework.

These contextual limitations highlight the importance of philosophical pluralism in approaching algorithmic fairness, with different frameworks being more or less appropriate depending on the specific context and the values at stake.

## 6. Conclusion

### 6.1 Summary of Key Contributions

This thesis has developed a framework for algorithmic fairness based on Robert Nozick's entitlement theory of justice, offering a libertarian perspective that complements the predominantly egalitarian approaches in the existing literature. The key contributions of this work include:

1. **Conceptual Framework**: The thesis has articulated a conceptual framework for understanding algorithmic fairness from a Nozickian perspective, translating the principles of just acquisition, transfer, and rectification into the context of algorithmic decision-making.

2. **Formal Metrics**: The thesis has proposed formal metrics for evaluating algorithmic fairness according to Nozickian principles, including Procedural Entitlement Fairness (PEF), Consent and Transparency Index (CTI), Rectification Responsiveness Measure (RRM), and Choice Sensitivity Ratio (CSR).

3. **Integrated Approach**: The thesis has developed an integrated Nozickian Fairness Score (NFS) that combines these metrics into a comprehensive

evaluation framework, allowing for context-specific weighting of different components.

4. **Case Studies**: The thesis has applied the Nozickian framework to hypothetical case studies in algorithmic lending and hiring, demonstrating its practical utility and highlighting contrasts with traditional egalitarian approaches.

5. **Philosophical Analysis**: The thesis has examined the philosophical implications of adopting a Nozickian perspective on algorithmic fairness, including implications for conceptions of procedural vs. distributive justice, individual rights vs. group fairness, and the role of history in evaluating algorithmic systems.

These contributions expand the philosophical foundations of algorithmic fairness research, offering a novel perspective that foregrounds procedural justice, individual rights, and historical context.

### 6.2 Broader Implications

The Nozickian framework developed in this thesis has several broader implications for the field of algorithmic fairness and ethics:

**Philosophical Pluralism**: By introducing a libertarian perspective into a field dominated by egalitarian approaches, this thesis highlights the importance of philosophical pluralism in algorithmic ethics. Different philosophical traditions offer distinct insights and values that can enrich our understanding of what constitutes fair algorithmic treatment.

**Balancing Process and Pattern**: The framework suggests the need to balance concerns about procedural justice (how decisions are made) with concerns about distributive patterns (what outcomes result). Both aspects are morally significant, and a comprehensive approach to algorithmic fairness should attend to both.

**Historical Context**: The framework emphasizes the importance of historical context in evaluating algorithmic systems, suggesting that we cannot assess fairness without considering how data was acquired, how past injustices might affect current decisions, and what rectification might be required.

**Individual and Group Perspectives**: The framework highlights tensions between individual and group perspectives on justice, suggesting that both have important roles in evaluating algorithmic systems but that they may sometimes point in different directions.

These broader implications suggest that algorithmic fairness research should embrace greater philosophical diversity, contextual sensitivity, and normative clarity about the values and principles that inform different approaches to fairness.

### 6.3 Future Research Directions

This thesis points to several promising directions for future research:

**Empirical Testing**: Future work could empirically test the proposed Nozickian metrics on real-world algorithmic systems, exploring their practical utility and limitations.

**Cross-Philosophical Integration**: Future research could explore how Nozickian and egalitarian approaches to algorithmic fairness might be integrated or balanced in specific contexts, drawing on the strengths of each perspective.

**Contextual Refinement**: The Nozickian framework could be refined and adapted for specific domains (e.g., healthcare, education, criminal justice), attending to the particular rights, entitlements, and historical injustices relevant to each context.

**Critique and Response**: Future work could engage more deeply with critiques of Nozick's theory and explore how the framework might be modified to address these critiques while maintaining its core insights.

**Regulatory Applications**: Research could examine how the Nozickian framework might inform regulatory approaches to algorithmic systems, balancing concerns about individual rights and procedural justice with other social values.

These research directions would build on the foundation laid in this thesis, further exploring the potential contributions of libertarian political philosophy to algorithmic ethics and fairness.

### 6.4 Concluding Reflections

The increasing role of algorithms in social decision-making raises profound questions about justice, fairness, and the proper relationship between individuals and automated systems. As we navigate these questions, we need philosophical frameworks that can help us articulate and balance the diverse values at stake.

This thesis has argued that Nozick's entitlement theory, despite its relative absence from current discussions of algorithmic fairness, offers valuable insights for this endeavor. By focusing on procedural justice, individual rights, and historical context, a Nozickian approach complements existing egalitarian perspectives, enriching our understanding of what constitutes fair algorithmic treatment.

As we continue to develop and deploy algorithmic systems that affect human lives, we should draw on the full range of philosophical traditions to guide our ethical evaluations and technical innovations. By incorporating diverse perspectives, including both egalitarian and libertarian approaches, we can work toward algorithmic systems that respect the complex and multifaceted nature of justice in the digital age.

# References

Angwin, J., Larson, J., Mattu, S., & Kirchner, L. (2016). Machine bias. ProPublica, 23(2016), 139-159.

Arneson, R. J. (1989). Equality and equal opportunity for welfare. Philosophical Studies, 56(1), 77-93.

Barocas, S., & Selbst, A. D. (2016). Big data's disparate impact. California Law Review, 104, 671-732.

Barocas, S., Hardt, M., & Narayanan, A. (2019). Fairness and machine learning. https://fairmlbook.org/

Binns, R. (2018). Fairness in machine learning: Lessons from political philosophy. In Conference on Fairness, Accountability and Transparency (pp. 149-159).

Binns, R. (2020). On the apparent conflict between individual and group fairness. In Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (pp. 514-524).

Bogen, M., & Rieke, A. (2018). Help wanted: An examination of hiring algorithms, equity, and bias. Upturn.

Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. Big Data, 5(2), 153-163.

Cohen, G. A. (1989). On the currency of egalitarian justice. Ethics, 99(4), 906-944.

Cohen, G. A. (1995). Self-ownership, freedom, and equality. Cambridge University Press.

Corbett-Davies, S., & Goel, S. (2018). The measure and mismeasure of fairness: A critical review of fair machine learning. arXiv preprint arXiv:1808.00023.

Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. In Proceedings of the 3rd innovations in theoretical computer science conference (pp. 214-226).

Dworkin, R. (2000). Sovereign virtue: The theory and practice of equality. Harvard University Press.

Eubanks, V. (2018). Automating inequality: How high-tech tools profile, police, and punish the poor. St. Martin's Press.

Fazelpour, S., & Lipton, Z. C. (2020). Algorithmic fairness from a non-ideal perspective. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (pp. 57-63).

Frankfurt, H. (1987). Equality as a moral ideal. Ethics, 98(1), 21-43.

Fuster, A., Goldsmith-Pinkham, P., Ramadorai, T., & Walther, A. (2022). Predictably unequal? The effects of machine learning on credit markets. The

Journal of Finance, 77(1), 5-47.

Green, B. (2018). "Fair" risk assessments: A precarious approach for criminal justice reform. In 5th Workshop on Fairness, Accountability, and Transparency in Machine Learning.

Grgić-Hlača, N., Redmiles, E. M., Gummadi, K. P., & Weller, A. (2018). The case for process fairness in learning: Feature selection for fair decision making. *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 51–57. https://doi.org/10.1145/3278721.3278725

Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems* (pp. 3315–3323).

Kleinberg, J., Mullainathan, S., & Raghavan, M. (2017). Inherent trade-offs in the fair determination of risk scores. In *Proceedings of the 8th Innovations in Theoretical Computer Science Conference* (ITCS).

Lee, M., & Floridi, L. (2021). Algorithmic fairness in AI for social good: A survey. *Philosophy & Technology*, 34, 1023–1052.

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.

Otsuka, M. (2003). *Libertarianism without inequality*. Oxford University Press.

Parfit, D. (1997). Equality and priority. *Ratio*, 10(3), 202–221.

Rawls, J. (1971). *A theory of justice*. Harvard University Press.

Selbst, A. D., Boyd, D., Friedler, S. A., Venkatasubramanian, S., & Vertesi, J. (2019). Fairness and abstraction in sociotechnical systems. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (pp. 59–68).

Thierer, A. (2016). Permissionless innovation: The continuing case for comprehensive technological freedom. *Mercatus Center at George Mason University.*

Waldron, J. (1992). Superseding historic injustice. *Ethics*, 103(1), 4–28.

Kymlicka, W. (2002). *Contemporary political philosophy: An introduction* (2nd ed.). Oxford University Press.

Nielsen, K. (1979). Radical egalitarian justice: Justice as equality. *Social Theory and Practice*, 5(2), 209–226.