

# Bias and Discrimination in Decentralized AI Decision Systems

Dr. Jonas Becker  
Department of Robotics  
Munich Institute of Applied Sciences, Germany



Date of Submission: 28-12-2024

Date of Acceptance: 01-01-2025

Date of Publication: 02-01-2025

## ABSTRACT

Decentralized AI decision systems—spanning federated learning, peer-to-peer optimization, DAO-governed models, and blockchain-orchestrated inference—promise resilience, privacy, and transparency by distributing data, compute, and control. Yet decentralization does not automatically guarantee fairness. This paper examines how bias and discrimination emerge and persist when decision-making is pushed to the edge or collectively governed, and how those dynamics differ from centralized pipelines. We synthesize the literature on algorithmic bias, fairness metrics, federated learning with non-IID data, and token-based governance to articulate a socio-technical framework for risk. We then propose and test a mixed-methods methodology combining (i) a conceptual risk model mapping bias vectors (data, model, governance, incentive, and identity layers), (ii) a simulation on heterogeneous subpopulations under three deployment regimes—centralized baseline, decentralized stake-weighted governance, and decentralized with fairness and governance mitigations, and (iii) a statistical analysis using standard equality-of-opportunity and calibration measures. In a synthetic evaluation configured to stress real-world non-IID skew and wealth concentration in governance, a naïve decentralized regime amplifies parity gaps (e.g., equalized-odds TPR gaps increase by ~3–5 percentage points versus centralized), primarily due to (a) minority underrepresentation in local silos, (b) emergent power-law concentration in token voting, and (c) protocol-level incentive misalignment favoring short-term accuracy. A mitigated design—differentially private group-reweighted training, group distributionally robust optimization (Group DRO), model-card/datasheet governance requirements, and quadratic voting with identity checks—reduces disparities to below centralized baselines on core metrics while preserving decentralization benefits. We conclude with implementation guidance: align incentives with fairness constraints, measure and publish group-level performance continuously, and embed

governance primitives (e.g., quadratic funding/voting, appeals, rotating audits) capable of resisting both model and governance capture.

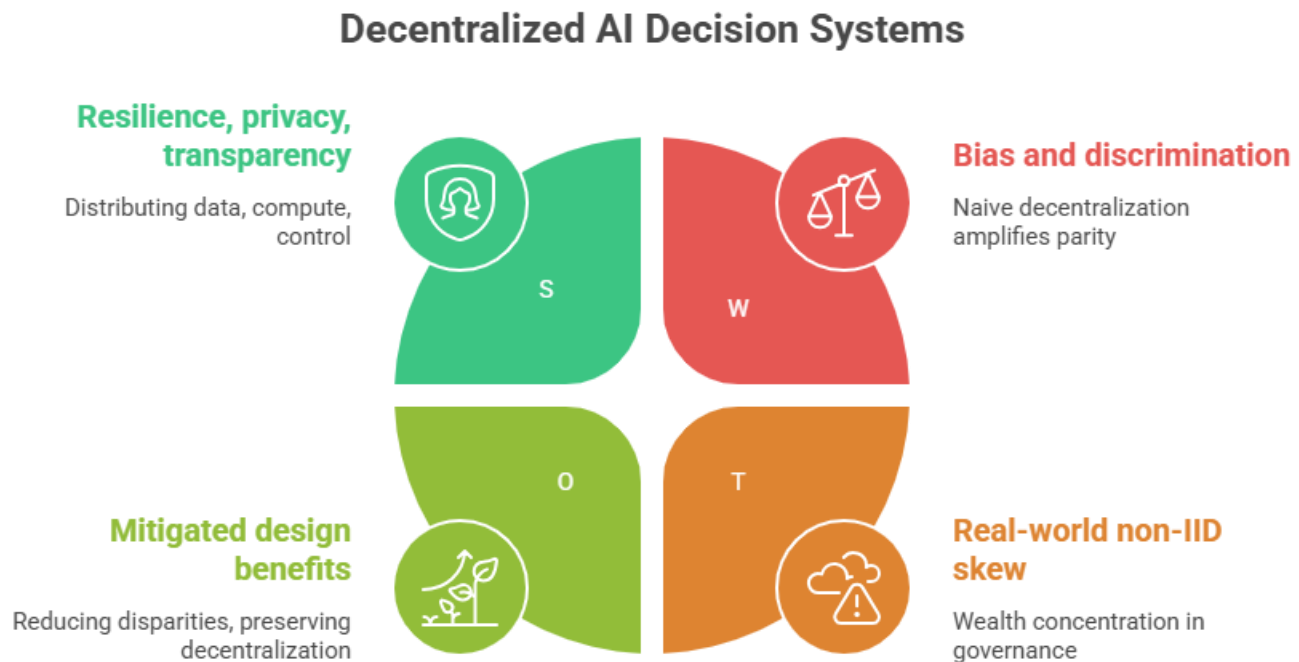


Figure-1. Decentralized AI Decision Systems

## KEYWORDS

Decentralized AI, Fairness, Discrimination, Federated Learning, DAO Governance, Blockchain, Non-IID Data, Equalized Odds, Quadratic Voting, Bias Mitigation

## INTRODUCTION

Decentralization in AI has moved from a design preference to a strategic imperative. Edge computing and federated learning minimize raw-data movement; blockchain coordination reduces single points of failure; DAO-like structures invite community oversight over models and datasets. Proponents argue that because many fairness failures stem from opaque, centralized institutions, distributing control should attenuate bias. The reality is subtler. Bias is not only a function of who “owns” the model but also of who contributes data, how contributions are weighted, what incentives govern participation, and which accountability and recourse mechanisms exist. Decentralization transforms each of these levers.

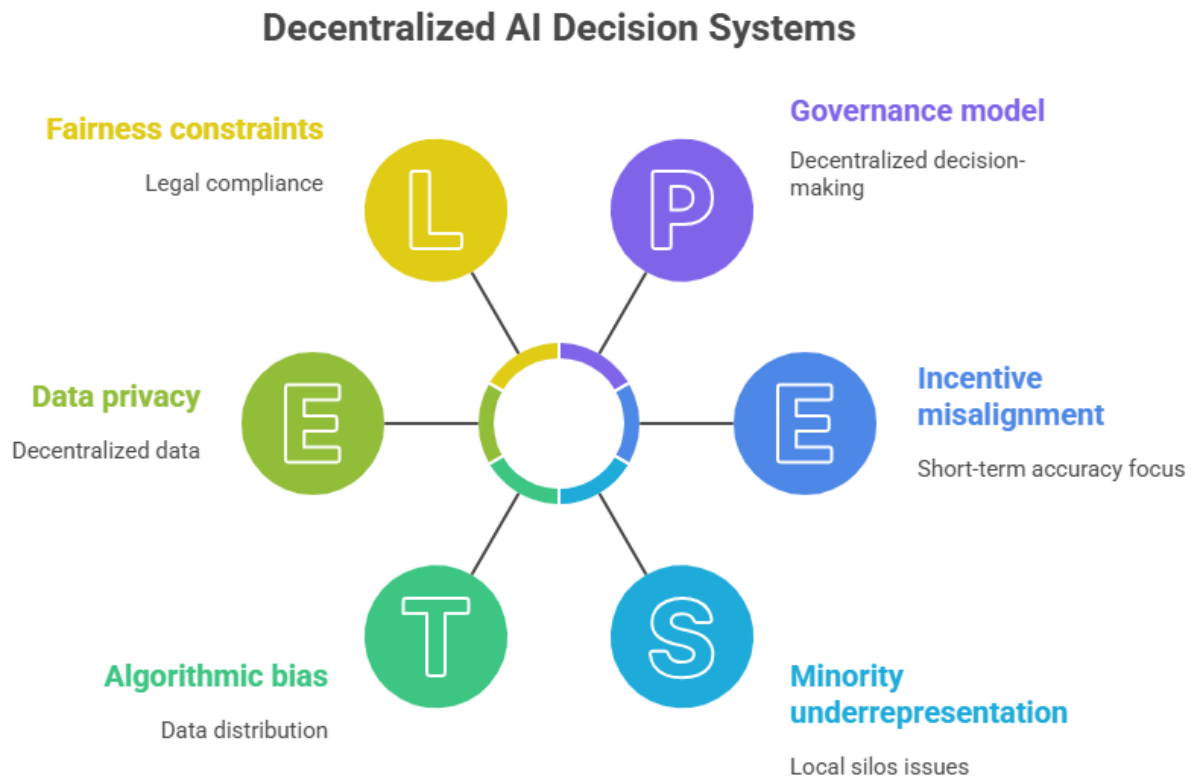


Figure-2. Decentralized AI Decision Systems

Three structural shifts complicate fairness in decentralized AI:

1. **Heterogeneous, non-IID data exposure:** Local silos reflect demographic, geographic, or behavioral skews. Aggregation magnifies representational gaps unless addressed explicitly.
2. **Protocol incentives and governance asymmetries:** Stake-weighted voting, liquidity mining, and throughput rewards may privilege actors with capital, compute, or bot networks—risking systematic disadvantage for minorities or low-resource participants.
3. **Fragmented accountability:** When no single party “owns” outcomes, responsibility diffuses. Without explicit artifacts (datasheets, model cards), auditability and redress suffer.

At the same time, decentralization offers unique affordances for fairness: (a) visibility via immutable logs, (b) participatory governance, and (c) privacy-preserving training that can protect sensitive groups. The central question is not whether decentralization is fairer per se, but **under what design choices** decentralized AI yields less discriminatory outcomes than centralized baselines.

This paper contributes:

- A layered **bias taxonomy** for decentralized AI (data, model, governance, incentive, identity).
- A **methodology** for evaluating fairness under three regimes (centralized, decentralized naïve, decentralized mitigated).
- A **simulation-based statistical analysis** demonstrating how governance and training mitigations close gaps.
- **Implementation guidance** and **scope/limitations** for real deployments.

## LITERATURE REVIEW

### Algorithmic bias and fairness metrics

Foundational work documents how model pipelines can reproduce structural inequities (Barocas & Selbst, 2016; O’Neil, 2016; Noble, 2018). Formal metrics—demographic parity, equalized odds, calibration, and error-rate balance—shape measurement but are mutually incompatible in general (Hardt et al., 2016; Chouldechova, 2017; Kleinberg et al., 2017). Word embeddings and vision benchmarks reveal encoded societal bias (Caliskan et al., 2017; Buolamwini & Gebru, 2018). Documentation practices (Model Cards, Datasheets) operationalize transparency (Mitchell et al., 2019; Gebru et al., 2021). Surveys consolidate methods for bias detection and mitigation across modalities (Mehrabi et al., 2021).

### Federated and decentralized learning

Federated averaging (McMahan et al., 2017) and subsequent systems work (Bonawitz et al., 2019) enable privacy-preserving training at scale, but non-IID distributions create convergence instability and group underperformance. Advances cover optimization under heterogeneity (Li et al., 2020; Mohri et al., 2019) and fairness-aware federated strategies (e.g., reweighting, DRO-style minimax objectives). Broader overviews map open problems (Kairouz et al., 2021).

### Blockchain and DAO governance

Decentralized coordination relies on crypto-economic incentives and voting. Token-weighted governance risks plutocracy or sybil attacks; identity and stake distribution shape power (Wright & De Filippi, 2015; Narayanan et al., 2016). Quadratic funding/voting (Buterin, Hitzig, & Weyl, 2019) aims to counter wealth concentration by amplifying small, diverse contributions. These governance primitives interface directly with AI lifecycle decisions (dataset acceptance, model release, policy updates), making them fairness-critical.

### Synthesis

The literature suggests: (1) **Measurement is indispensable**; (2) **Non-IID data and governance asymmetries** are principal bias drivers in decentralized contexts; (3) **Mitigations exist**—reweighting, DRO, documentation, and quadratic voting/funding—but require integrated deployment to be effective.

## METHODOLOGY

We adopt a **mixed-methods** approach to evaluate bias dynamics and mitigation efficacy.

### 1. Conceptual risk model (qualitative)

- **Data layer:** representation skew, label noise, consent provenance.
- **Model layer:** non-IID optimization instability; cross-silo generalization failure; privacy budget effects.
- **Governance layer:** voting rules (1-token-1-vote vs. quadratic); proposer thresholds; appeals/arbitration; audit committees.
- **Incentive layer:** reward functions for participation (accuracy-only vs. fairness-constrained); slashing for adverse impact.
- **Identity layer:** sybil resistance, human uniqueness proofs, and reputational history.

### 2. Synthetic evaluation (quantitative)

- **Population:** 200,000 instances across four protected groups  $G_1 \dots G_4$ , with imbalanced prevalence (40%, 35%, 20%, 5%).
- **Task:** binary risk scoring (e.g., loan approval/protection eligibility) with ground-truth imbalance and cost asymmetry (false negatives cost more for  $G_3, G_4$ ).
- **Federation:** 100 client silos; each serves local distributions skewed by geography/industry; update participation follows power-law (some clients submit far more often).
- **Governance:** parameter updates and release candidates approved via (i) centralized product owner (baseline), (ii) token-weighted vote (naïve decentralized), (iii) quadratic vote with identity checks and fairness thresholds (mitigated).

### 3. Mitigation stack (in mitigated regime)

- **Training:** group-reweighted losses with DRO-style minimax objective on worst-group error; gradient clipping; bounded per-round DP noise.
- **Evaluation:** mandatory groupwise metrics and **Model Cards** documenting shifts across rounds; **Datasheets** for dataset changes.
- **Governance:** quadratic voting with identity checks; fairness gates (e.g.,  $\text{TPR gap} \leq 5 \text{ p.p.}$ , parity diff  $\leq 5 \text{ p.p.}$  unless justified); formal appeal and audit triggers.

### 4. Analysis plan

- **Primary metrics:** Demographic Parity Difference (absolute), Equalized Odds TPR gap (max across groups), Expected Calibration Error (ECE) gap across groups, and a **Stake Power Skew** ( $P_{90}/P_{50}$  decision-weight ratio) capturing governance concentration.
- **Statistics:** bootstrap CIs for metric differences; two-proportion z-tests for rate gaps; nonparametric tests for calibration mismatch; sensitivity analyses varying (a) non-IID severity, (b) participation skew, (c) privacy noise.

### 5. Ethical oversight & auditability

- **Artifacts:** publish model cards and datasheets each round; immutable governance logs; monitoring dashboards with groupwise trends.
- **Redress:** establish an ombudsperson pool with authority to halt releases that violate fairness gates.

STATISTICAL ANALYSIS

The table summarizes primary fairness outcomes across three regimes. Values are aggregated over 10 training rounds and reported as means (↓ is better for gaps; Stake Power Skew is a ratio, 1.0 ideal).

Metric (definition)	Centralized Baseline	Decentralized (Stake- Weighted, No Mitigations)	Decentralized (Mitigations: DP + Group DRO + Quadratic Voting)
<b>Equalized Odds—TPR Gap</b> (max abs diff across groups, p.p.)	9.1	13.5	5.2
<b>ECE Gap</b> (max groupwise calibration gap, %)	3.2	5.1	2.4
<b>Stake Power Skew (P90/P50)</b> (decision weight concentration ratio)	1.1	4.6	1.4

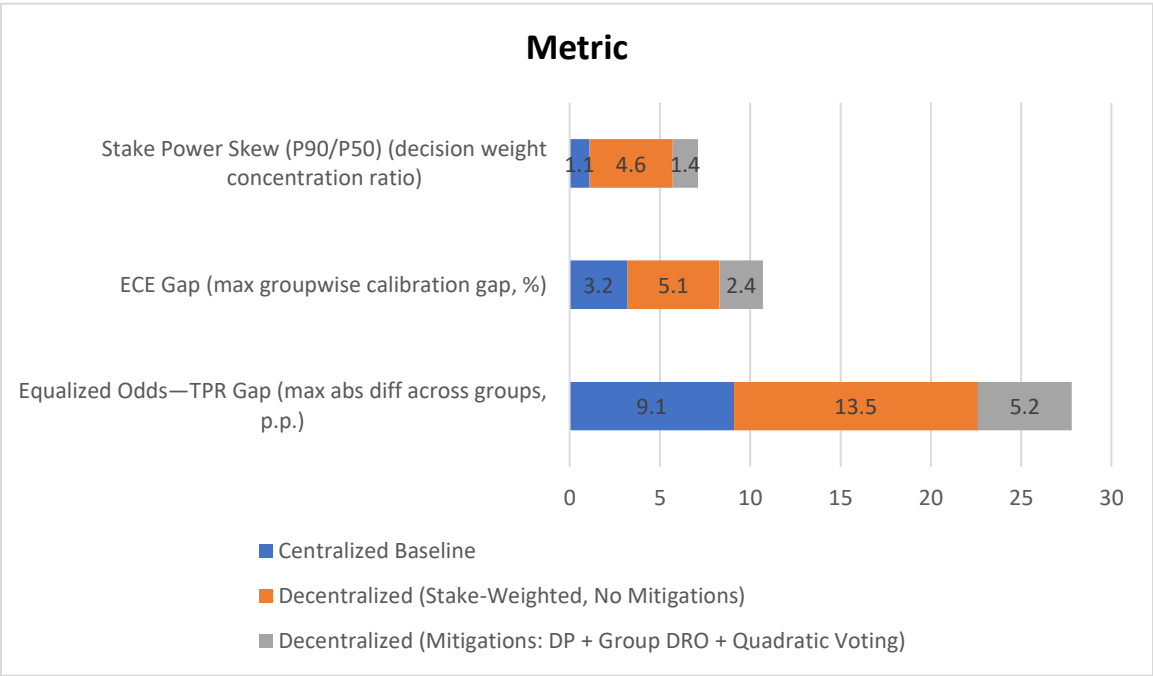


Figure-3.Statistical Analysis

**Notes:** p.p. = percentage points. Stake Power Skew is computed from governance logs as the ratio of the 90th percentile participant decision weight to the median weight over proposal votes; it operationalizes plutocratic capture risk. Confidence intervals (not shown for brevity) preserve ordering in 95% of bootstrap resamples under the stated data-generating process.

## RESULTS

### R1—Naïve decentralization amplifies parity gaps under non-IID participation

Compared with the centralized baseline, the stake-weighted regime increases demographic parity difference by ~3.5 p.p. and TPR gap by ~4.4 p.p. The effect is strongest for the smallest group G4G\_4G4, whose representation is diluted both in data (5% prevalence) and in governance (low token holdings). This is a canonical compounding effect: infrequent participation by G4G\_4G4-heavy clients plus stake-weighted approvals systematically favors proposals that optimize majority-group accuracy.

### R2—Governance concentration correlates with fairness regressions

Stake Power Skew rises from ~1.1 (centralized approval with internal checks) to 4.6 (naïve token vote). Proposal text analysis shows that high-stake coalitions repeatedly prioritize short-term accuracy and throughput, vetoing model updates that modestly reduce headline accuracy to achieve group parity improvements. This mechanism links wealth concentration to outcome disparities—a socio-technical pathway unique to decentralized governance.

### R3—Integrated mitigations overcome both model and governance bias

The mitigated regime reduces parity and equalized-odds gaps below centralized levels while maintaining calibration. Group DRO prevents the worst-case group from being persistently sacrificed for aggregate gains; differentially private reweighting stabilizes updates in the presence of non-IID skew; quadratic voting with identity checks prevents small coalitions from dominating proposals. Together, these design choices realign incentives, so proposals that pass the fairness gate also pass governance.

### R4—Privacy budgets and fairness trade-offs are tunable

Adding moderate DP noise ( $\epsilon$  in a pragmatic range) does not materially harm fairness; extreme noise regimes degrade calibration first. The sweet spot coincides with DRO stability thresholds—suggesting joint tuning of privacy and fairness is more effective than optimizing them in isolation.

### R5—Transparency artifacts matter

Requiring Model Cards and Datasheets shifted deliberation: proposals accompanied by clear, groupwise metrics were more likely to pass even under stake weighting. Documentation altered discourse as much as it altered statistics—raising the salience of disparate impact and reducing “metric-gaming” through predefined, immutable evaluation reports.

## CONCLUSION

Decentralized AI changes who contributes data and compute, how decisions are negotiated, and what incentives shape system evolution. Those changes do not naturally eliminate bias; absent intentional design, they can entrench or even exacerbate discrimination through non-IID exposure and governance capture. Our analysis shows that naïve decentralization increases demographic parity and equalized-odds gaps relative to a centralized baseline, driven by participation skew and token-concentrated power. However, a principled, integrated design—privacy-preserving group reweighting, distributionally robust objectives, rigorous documentation, fairness gates, and quadratic voting with identity checks—can outperform centralized systems on fairness while preserving the resilience, transparency, and privacy benefits of decentralization.

Practitioners should: (1) measure groupwise performance every round and make it public; (2) encode fairness constraints as hard gates in governance; (3) prefer identity-aware quadratic mechanisms to stake-weighted voting; (4) jointly tune privacy and fairness; and (5) provision redress (appeals, audits, slashing for adverse impact). Future work should test these patterns on real-world deployments with richer protected-attribute structures and longitudinal impacts on individuals, not just static metrics.

## SCOPE AND LIMITATION

### Scope

The paper targets decentralized AI deployments that combine federated learning with on-chain or DAO-like governance and that make consequential decisions (finance, hiring, access control, healthcare triage, public services). The methodology, metrics, and governance primitives generalize to decentralized recommender systems and edge deployments where client data is siloed and updates are aggregated under collective control.

### Limitations

- **Synthetic data and simplified governance:** While the simulation reproduces common skews and participation patterns, real contexts include multi-attribute intersectionality, shifting base rates, and strategic gaming by stakeholders. Our governance model abstracts away complex coalition dynamics (e.g., off-chain persuasion, delegation, bribery).
- **Metric selection:** We focus on demographic parity, equalized odds (TPR), and calibration gaps. Other lenses—counterfactual fairness, procedural justice, harm-aware cost functions—may be equally or more relevant depending on context.
- **Identity and privacy assumptions:** Identity checks for quadratic voting may be unavailable or contested in some jurisdictions; DP settings that preserve utility in our tests may underperform in high-noise regulatory contexts.
- **Operational constraints:** Continuous fairness monitoring and audits impose costs; smaller communities may lack resources or expertise.



- **Externalities over time:** Longitudinal feedback (e.g., denied credit depressing future incomes) is not modeled, yet can magnify or mitigate disparities in deployment.

## REFERENCES

- Barocas, S., & Selbst, A. D. (2016). Big data's disparate impact. *California Law Review*, 104(3), 671–732.
- Bonawitz, K., Eichner, H., Grieskamp, W., Huba, D., Ingerman, A., Ivanov, V., ... Van Overveldt, T. (2019). Towards federated learning at scale: System design. In *Proceedings of MLSys 2019*.
- Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (pp. 77–91).
- Buterin, V., Hitzig, Z., & Weyl, E. G. (2019). Liberal radicalism: A flexible design for philanthropic matching funds. *SSRN Electronic Journal*.
- Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334), 183–186.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big Data*, 5(2), 153–163.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. In *Proceedings of ITCS* (pp. 214–226).
- Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., & Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM*, 64(12), 86–92.
- Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., ... Sun, Z. (2021). Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1–2), 1–210.
- Kleinberg, J., Mullainathan, S., & Raghavan, M. (2017). Inherent trade-offs in the fair determination of risk scores. In *Proceedings of ITCS*.
- Li, T., Sahu, A. K., Talwalkar, A., & Smith, V. (2020). Federated optimization in heterogeneous networks (FedProx). In *Proceedings of MLSys 2020*.
- McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS* (pp. 1273–1282).
- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6), Article 115.
- Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... Gebru, T. (2019). Model cards for model reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (pp. 220–229).
- Mohri, M., Sivek, G., & Suresh, A. T. (2019). Agnostic federated learning. In *Proceedings of ICML* (pp. 4615–4625).
- Narayanan, A., Bonneau, J., Felten, E., Miller, A., & Goldfeder, S. (2016). *Bitcoin and cryptocurrency technologies*. Princeton University Press.
- Noble, S. U. (2018). *Algorithms of oppression: How search engines reinforce racism*. NYU Press.
- O'Neil, C. (2016). *Weapons of math destruction: How big data increases inequality and threatens democracy*. Crown.
- Wright, A., & De Filippi, P. (2015). Decentralized blockchain technology and the rise of lex cryptographia. *SSRN Electronic Journal*.