

# Consent-Theoretic Framework for Quantifying Legitimacy

*Stakes, Voice, and Friction in Adversarial Governance*

Murad Farzulla<sup>1,2,\*</sup>

<sup>1</sup>Dissensus AI, London, UK    <sup>2</sup>King's College London, London, UK

\*Correspondence: [murad@dissensus.ai](mailto:murad@dissensus.ai)    ORCID: 0009-0002-7164-8704

February 2026

## Abstract

This paper develops a unified analytical framework for measuring political legitimacy across heterogeneous governance domains. Building on insights from constitutional political economy, social choice theory, and institutional analysis, the framework establishes consent-holding—the mapping from decision domains to those with authority over them—as a structural necessity of collective action. We formalize this intuition through five axioms and five theorems, demonstrating that legitimacy can be operationalized as stakes-weighted consent alignment  $\alpha(d, t)$ , while friction  $F(d, t)$  measures the deviation between outcomes and stakeholder preferences. The framework bridges normative democratic theory and empirical prediction, generating testable hypotheses about institutional stability. Historical validation examines suffrage expansion, abolition movements, labor rights, and contemporary platform governance, demonstrating how misalignment between stakes and voice generates observable instability. Unlike existing approaches that prescribe ideal institutions, this framework provides analytical tools for measuring legitimacy within any governance structure, enabling systematic comparison across democratic, technocratic, and algorithmic systems. Computational mechanism comparison via Bayesian learning dynamics across 1000 Monte Carlo runs demonstrates relative performance under adaptive agents: when preferences update based on observed policy outcomes, stakes-weighted DoCS achieves highest final alignment ( $\alpha = 0.872$ ) with lowest terminal friction ( $F = 1.5$ , 94.9% reduction from initial  $F = 30.3$ ). This comparative advantage holds across static baseline ( $\alpha = 0.627$ ), learning dynamics ( $\alpha = 0.872$ ), and alternative temporal mechanisms, suggesting stakes-weighting produces superior initial matches that persist even when agents adapt to institutional performance. The framework's domain-specific approach resolves the apparent tension between consent and competence, showing both as complementary dimensions of institutional legitimacy. This framework is part of the Adversarial Systems Research program, which examines stability, alignment, and friction dynamics in complex systems where competing interests generate structural conflict.

**Keywords:** legitimacy, consent, political stability, social choice, institutional design, friction, stakes-weighting

**JEL Codes:** D70 (Social Choice), D71 (Social Choice; Clubs; Committees; Associations), P16 (Political Economy)

## Research Context

This work forms part of the Adversarial Systems Research program, which investigates stability, alignment, and friction dynamics in complex systems where competing interests generate structural conflict. The program examines how agents with divergent preferences interact within institutional constraints across multiple domains: political governance (this paper), financial markets (cryptocurrency volatility and regulatory responses), human cognitive development (trauma as maladaptive learning from adversarial training environments), and artificial intelligence alignment (multi-agent systems with competing objectives).

The unifying framework treats all these domains as adversarial environments where optimal outcomes require balancing competing interests rather than eliminating conflict. In political systems, this manifests as the tension between stakeholder consent and technocratic competence. In financial markets, it appears as the conflict between regulatory stability and market innovation. In human development, it emerges as the challenge of learning accurate models from noisy or adversarial training data. In AI systems, it surfaces as the alignment problem when multiple agents optimize for different reward functions.

The Doctrine of Consensual Sovereignty presented here provides the theoretical foundation for analyzing legitimacy in any adversarial environment by formalizing the relationship between stakes, voice, and friction. Future work will extend this framework to algorithmic governance systems, multi-stakeholder climate negotiations, and autonomous agent coordination problems where consent structures remain undefined but friction dynamics are already observable.

## Key Notation

| Symbol        | Definition   |
|---------------|--|
| $H_t(d)$      | Consent-holder mapping (who decides in domain $d$ at time $t$ )                    |
| $s_i(d)$      | Stakes of agent $i$ in domain $d$ (material/capability exposure)                   |
| $C_{i,d}$     | Consent power of agent $i$ in domain $d$ (decision authority)                      |
| $\alpha(d,t)$ | Consent alignment (stakes-weighted share of voice held by stakeholders)            |
| $F(d,t)$      | Friction (stakes-weighted deviation between outcomes and preferences)              |
| $L(d,t)$      | Legitimacy ( $w_1 \cdot \alpha + w_2 \cdot P$ , balancing consent and performance) |
| $P(d,t)$      | Performance/competence metric (domain-specific outcome quality)                    |
| $x_{i,d}^*$   | Agent $i$ 's ideal action in domain $d$ (preference)                               |
| $x_d(t)$      | Realized action/outcome in domain $d$ at time $t$                                  |

## Scope and Limitations

This paper presents a conceptual framework for analyzing legitimacy in adversarial environments. While we formalize core relationships mathematically—consent alignment  $\alpha(d,t)$ , friction  $F(d,t)$ , and legitimacy  $L(d,t)$ —complete operational measurement remains ongoing empirical work. Our contribution is providing analytical architecture that makes legitimacy comparable across governance domains, enabling systematic analysis previously confined to domain-specific theories.

We demonstrate proof-of-concept through computational validation via agent-based simulation and qualitative validation across seven historical domains (suffrage, abolition, labor rights, civil rights, LGBT inclusion, platform governance, climate policy). **Methodological note:** The computational models assume agents learn from outcomes (Bayesian updating), which by construction reduces friction as preferences converge—this compares which consent mechanisms produce superior alignment given plausible behavioral assumptions, not whether friction *can* reduce (that follows definitionally). Measurement challenges are acknowledged in Section 4.

This v1.0.0 preprint establishes theoretical foundations and demonstrates implementability; subsequent versions will expand empirical validation with fully quantified historical cases, refine measurement protocols, and extend applications to algorithmic governance, climate negotiations, and multi-agent AI systems.

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>7</b>  |
| <b>2</b> | <b>Literature Review and Theoretical Foundations</b>                      | <b>8</b>  |
| 2.1      | Constitutional Political Economy . . . . .                                | 9         |
| 2.2      | Social Choice Theory and Impossibility Results . . . . .                  | 9         |
| 2.3      | Stakeholder Theory and Corporate Governance . . . . .                     | 10        |
| 2.4      | Common-Pool Resource Governance . . . . .                                 | 11        |
| 2.5      | Deliberative Democracy and Mini-Publics . . . . .                         | 11        |
| 2.6      | Algorithmic Governance and Platform Legitimacy . . . . .                  | 12        |
| 2.7      | Voting Power Indices and Coalition Analysis . . . . .                     | 13        |
| 2.8      | Relational Autonomy and Consent Capacity . . . . .                        | 13        |
| <b>3</b> | <b>Formal Framework: Primitives, Axioms, and Theorems</b>                 | <b>14</b> |
| 3.1      | Primitives and Definitions . . . . .                                      | 14        |
| 3.2      | Axioms . . . . .  | 15        |
| 3.3      | Theorem 1: Consent-Holding Necessity . . . . .                            | 15        |
| 3.4      | Theorem 2: Inevitable Friction . . . . .                                  | 15        |
| 3.5      | Definition 1: Legitimacy as Consent Alignment . . . . .                   | 16        |
| 3.6      | Postulate 1: Competence-Consent Trade-Off . . . . .                       | 16        |
| 3.7      | Theorem 3: Minimal Absolutism from Relativism . . . . .                   | 17        |
| <b>4</b> | <b>Operationalization: Empirical Measurement and Identification</b>       | <b>17</b> |
| 4.1      | Formal Measurement Framework . . . . .                                    | 17        |
| 4.2      | Friction Metrics and Tolerance-Weighted Extensions . . . . .              | 17        |
| 4.3      | Empirical Identification Strategies . . . . .                             | 18        |
| 4.4      | Testable Predictions and Empirical Hypotheses . . . . .                   | 18        |
| <b>5</b> | <b>Social Contract Theories as Distribution Mechanisms</b>                | <b>18</b> |
| 5.1      | Rawlsian Justice as Maximin Consent . . . . .                             | 19        |
| 5.2      | Utilitarian Consent as Weighted Aggregation . . . . .                     | 19        |
| 5.3      | Libertarian Consent as Property Rights . . . . .                          | 19        |
| <b>6</b> | <b>Historical Validation: Case Studies in Consent Alignment Dynamics</b>  | <b>19</b> |
| 6.1      | Suffrage Expansion: Gradual Consent Broadening . . . . .                  | 19        |
| 6.2      | Abolition Movements: Maximum Stakes, Zero Consent . . . . .               | 20        |
| 6.3      | Labor Rights and Corporate Codetermination . . . . .                      | 20        |
| 6.4      | Platform Governance Rebellions (2010s-Present) . . . . .                  | 21        |
| 6.5      | Scope Conditions: When Friction Fails to Generate Incorporation . . . . . | 21        |
| <b>7</b> | <b>Computational Mechanism Comparison: Adaptive Learning Dynamics</b>     | <b>23</b> |
| 7.1      | Simulation Design . . . . .   | 24        |
| 7.2      | Bayesian Preference Learning Dynamics . . . . .                           | 24        |
| 7.3      | Results . . . . .   | 24        |
| 7.3.1    | Static Baseline Comparison . . . . .                                      | 24        |

|           |  |           |
|-----------|--|-----------|
| 7.3.2     | Bayesian Learning Dynamics: Genuine Convergence . . . . .                | 25        |
| <b>8</b>  | <b>Dynamic Validation and Robustness</b>                                 | <b>26</b> |
| 8.1       | Convergence Statistics . . . . .   | 26        |
| 8.2       | Robustness Across Dynamic Mechanisms . . . . .                           | 27        |
| 8.3       | Plutocracy Convergence: Co-option Versus Legitimacy . . . . .            | 28        |
| 8.4       | Robustness to Parameter Variations . . . . .                             | 28        |
| <b>9</b>  | <b>Objections and Replies</b>  | <b>29</b> |
| 9.1       | Objection 1: Infinite Regress . . . . .                                  | 29        |
| 9.2       | Objection 2: Stakes Manipulation (Plutocracy) . . . . .                  | 29        |
| 9.3       | Objection 3: Competence Sacrifice . . . . .                              | 29        |
| 9.4       | Objection 4: Unresponsive Minorities . . . . .                           | 29        |
| 9.5       | Objection 5: Future Generations . . . . .                                | 29        |
| 9.6       | Objection 6: Collective Action Problems . . . . .                        | 29        |
| 9.7       | Objection 7: Cultural Relativism . . . . .                               | 30        |
| <b>10</b> | <b>Conclusion</b>  | <b>30</b> |
| 10.1      | Weight Determination as Endogenous Constitutional Problem . . . . .      | 31        |
| <b>A</b>  | <b>Appendix A: Robustness Checks</b>                                     | <b>32</b> |
| <b>B</b>  | <b>Appendix B: Extended Literature Synthesis</b>                         | <b>33</b> |
| B.1       | B.1 Democratization and Institutional Sequencing . . . . .               | 33        |
| B.2       | B.2 Protest Dynamics and Threshold Mobilization . . . . .                | 33        |
| B.3       | B.3 Deliberative Capacity and Epistemic Inclusion . . . . .              | 33        |
| B.4       | B.4 Algorithmic Governance, Opacity, and Contestation . . . . .          | 34        |
| B.5       | B.5 Climate and Intergenerational Governance . . . . .                   | 34        |
| B.6       | B.6 Summary for Journal Reduction . . . . .                              | 34        |
| <b>C</b>  | <b>Appendix C: Extended Dynamic Comparison Tables</b>                    | <b>34</b> |
| <b>D</b>  | <b>Appendix D: Methodological Claim Boundaries</b>                       | <b>35</b> |
| D.1       | D.1 Static Versus Dynamic Identification . . . . .                       | 35        |
| D.2       | D.2 Claims Defensible from Current Evidence . . . . .                    | 35        |
| D.3       | D.3 Claims Requiring Additional Design Before Strong Inference . . . . . | 35        |
| D.4       | D.4 Extensions Prioritized for Journal Version . . . . .                 | 36        |
| <b>E</b>  | <b>Appendix E: Additional Historical Exploration Cases</b>               | <b>36</b> |
| E.1       | E.1 Civil Rights Incorporation Dynamics . . . . .                        | 36        |
| E.2       | E.2 Climate Governance as Persistent Proxy Challenge . . . . .           | 36        |
| E.3       | E.3 Federal and Polycentric Designs . . . . .                            | 36        |
| E.4       | E.4 Exploratory Empirical Agenda . . . . .                               | 36        |
| <b>F</b>  | <b>Appendix F: Extended Social Contract Architecture</b>                 | <b>36</b> |
| F.1       | F.1 Four-Layer Interpretation Framework . . . . .                        | 37        |
| F.2       | F.2 Hobbesian Monopoly and Security-First Legitimacy . . . . .           | 37        |

|          |   |           |
|----------|---|-----------|
| F.3      | F.3 Lockean Conditional Delegation . . . . .                            | 37        |
| F.4      | F.4 Rousseauian General Will and Collective Self-Rule . . . . .         | 37        |
| F.5      | F.5 Technocratic Delegation as Expertise Concentration . . . . .        | 37        |
| F.6      | F.6 Anarchist and Federal Variants as Domain Fragmentation . . . . .    | 38        |
| F.7      | F.7 Algorithmic Social Contract and Code-Mediated Authority . . . . .   | 38        |
| F.8      | F.8 Comparative Matrix . . . . .  | 38        |
| F.9      | F.9 Endogenous Weight Selection Across Doctrines . . . . .              | 38        |
| F.10     | F.10 Practical Implication for Canonical and Journal Versions . . . . . | 38        |
| <b>G</b> | <b>Appendix G: Hobbes–Locke–Rousseau Text-to-Formal Mapping</b>         | <b>38</b> |
| G.1      | G.1 Mapping Template . . . . .  | 39        |
| G.2      | G.2 Hobbes (Order-First Authorization) . . . . .                        | 39        |
| G.3      | G.3 Locke (Conditional and Revocable Delegation) . . . . .              | 40        |
| G.4      | G.4 Rousseau (Collective Self-Rule) . . . . .                           | 40        |
| G.5      | G.5 Cross-Doctrine Comparative Reading . . . . .                        | 40        |
| G.6      | G.6 Minimal Empirical Coding Scheme . . . . .                           | 40        |

## 1 Introduction

Political legitimacy presents a fundamental puzzle: how can we measure whether authority is rightfully held across radically different governance domains? A state legislature, corporate board, algorithmic content moderation system, and common-pool resource management regime all make consequential decisions affecting stakeholders, yet existing frameworks struggle to provide unified analytical tools for assessing their legitimacy. Democratic theory emphasizes popular sovereignty, grounding legitimacy in consent of the governed (Locke, 1980; Rousseau, 1997; Rawls, 1971; Habermas, 1984), public choice highlights constitutional constraints (Buchanan and Tullock, 1962), while recent work on algorithmic governance introduces new challenges to consent-based legitimacy (Grimmelikhuijsen et al., 2022). What remains elusive is a framework capable of both normative evaluation and empirical prediction that applies consistently across domains.

This paper addresses this gap by developing consent-holding theory, an axiomatic framework that treats legitimacy as a structural property of decision-making systems rather than a binary classification. The central insight is deceptively simple yet powerful: in any domain where collective decisions produce shared consequences, someone must hold the authority to decide. This consent-holder mapping  $H_t(d)$ —identifying who decides in domain  $d$  at time  $t$ —is not a normative choice but a logical necessity arising from the structure of collective action itself. The framework’s contribution lies not in prescribing who should hold consent, but in providing rigorous tools for measuring the consequences of any particular allocation.

The framework makes three distinct contributions to political theory and institutional analysis. First, it establishes a formal connection between consent alignment and observable political friction. While democratic theorists have long argued that excluding affected stakeholders undermines legitimacy (Estlund, 2008), existing accounts lack operational metrics for testing these claims. We define consent alignment  $\alpha(d, t)$  as the stakes-weighted share of decision power held by affected parties, and friction  $F(d, t)$  as the stakes-weighted deviation between outcomes and stakeholder preferences. The framework predicts that persistent misalignment generates measurable instability—protests, non-compliance, institutional breakdown—making legitimacy empirically falsifiable rather than purely philosophical.

Second, the framework resolves the apparent tension between consent and competence through a competence-consent trade-off theorem (T4). Epistemic democrats argue that inclusive decision-making produces better outcomes through cognitive diversity (Landemore, 2013; Hong and Page, 2004), while critics worry that expanding consent sacrifices technical expertise. Our framework shows these concerns reflect different positions on the legitimacy frontier: some domains optimally weight performance highly (nuclear safety, pandemic response), while others prioritize consent alignment (constitutional amendments, community norms). Rather than declaring one approach universally superior, the framework provides tools for identifying domain-appropriate balances.

Third, this approach enables systematic historical and comparative analysis. By operationalizing legitimacy as  $\alpha(d, t)$  and friction as  $F(d, t)$ , we can trace institutional evolution quantitatively. Franchise expansion emerges not as discrete events but as gradual increases in  $\alpha(d)$  driven by the accumulating friction  $F(d)$  from excluding high-stakes populations. Women’s suffrage movements, abolition struggles, labor organizing, and contemporary platform governance rebellions all exhibit the same underlying dynamic: groups with high stakes  $s_i(d)$  but zero consent power  $C_i$  generate sustained friction until incorporation or suppression occurs. This pattern, predicted by the framework’s core theorems, provides empirical validation across centuries and continents.

The framework proceeds from seven minimal axioms to five core results establishing structural ne-

cessities (Section 3). Theorem 1 demonstrates consent-holding necessity: wherever decisions occur, some mapping  $H_i(d)$  must exist. Theorem 2 establishes inevitable friction: plural preferences guarantee that someone’s interests will be compromised unless perfect alignment obtains. Definition 1 operationalizes legitimacy as stakes-weighted consent alignment, providing an empirical metric. Postulate 1 formalizes the competence-consent trade-off, showing legitimacy as a weighted combination  $L = w_1 \cdot \alpha + w_2 \cdot P$ . Theorem 3 derives a minimal absolutism from value relativism: even if content-level values are frame-dependent, the existence of consent-holding structures remains invariant.

Section 2 situates this framework within nine research traditions—constitutional political economy (Buchanan and Tullock, 1962), social choice theory (Arrow, 1951; Sen, 2017), stakeholder theory (Freeman, 1984), common-pool resource governance (Ostrom, 1990), deliberative democracy (Fishkin, 2018; Habermas, 1990), algorithmic governance (Barocas et al., 2019), epistemic democracy (Estlund, 2008; Landemore, 2013), relational autonomy (Mackenzie, 2014), and legitimacy theory (Scharpf, 1999; Schmidt, 2013). Rather than competing with these approaches, consent-holding theory provides a unifying analytical architecture: each tradition contributes insights about how consent should be allocated or what constitutes legitimate use of authority, while our framework offers measurement tools applicable regardless of normative commitments.

Section 4 operationalizes the framework for empirical application, specifying proxy variables for consent power  $C_i$  (voting weights, agenda control, board representation), stakes  $s_i(d)$  (exposure measures, capability impacts, revealed preferences), and friction  $F(d)$  (protest incidence, litigation rates, policy reversals). This operationalization enables econometric identification strategies using panel data with institutional variation as instruments for consent alignment. Historical validation (Section 6) examines six cases spanning two centuries: women’s suffrage (1890s-1970s), abolition (1780s-1860s), labor rights (1850s-1930s), civil rights (1950s-present), LGBT inclusion (1969-present), and platform governance (2010s-present). Each case demonstrates the predicted pattern: high  $s_i(d)$  combined with zero  $C_i$  generates rising  $F(d)$  until elites respond through suppression or incorporation.

Section 9 addresses seven major objections, from concerns about infinite regress in consent structures to worries that stakes-weighting enables plutocracy. Section 10 concludes by outlining the research agenda this framework enables: cross-national legitimacy indices, institutional experiments varying  $\alpha(d)$  systematically, and applications to emerging domains (AI governance, climate policy, platform regulation) where consent structures remain contested.

The framework’s title—“consent-holding” rather than “consent theory”—reflects its analytical focus. This is not another account of why consent matters normatively, but a systematic investigation of how consent operates structurally. Just as markets emerge from property rights and contracts regardless of normative justifications for capitalism, consent-holding structures emerge from the necessity of collective decision-making regardless of democratic commitments. The framework’s power lies in making these structures visible, measurable, and comparable, enabling rigorous analysis of legitimacy claims that have historically remained philosophically contested but empirically elusive.

## 2 Literature Review and Theoretical Foundations

The consent-holding framework synthesizes and extends insights from nine distinct research traditions. This section reviews each tradition’s core contributions, identifies limitations the framework addresses, and demonstrates how operationalizing legitimacy as  $\alpha(d, t)$  and friction as  $F(d, t)$  enables empirical validation of long-standing theoretical claims.



## 2.1 Constitutional Political Economy

Building on earlier social contract foundations from [Hobbes \(1651\)](#) who established consent as prerequisite for legitimate political authority, Buchanan and Tullock’s [\(1962\)](#) seminal work establishes constitutional choice as a distinct analytical problem requiring different decision rules than ordinary politics. [Brennan and Buchanan \(1985\)](#) further develop this constitutional economics framework, distinguishing levels of collective action and establishing how constitutional rules create frameworks for collective choice. Their framework rests on several foundational insights that anticipate the consent-holding approach. First, they distinguish between constitutional rules—rarely changed frameworks establishing decision procedures—and political decisions made within those rules. This maps directly onto our concept of nested consent-holding:  $H_t(d_{meta})$  represents the consent-holders for constitutional domains, while  $H_t(d)$  operates within constraints established at the meta-level.

Second, Buchanan and Tullock argue that rational agents behind a “veil of uncertainty” would unanimously consent to rules benefiting all. Once constitutional structures are established, majority rule becomes acceptable for routine decisions. This anticipates our Theorem 1: consent-holding exists at every level, from object-level policy to constitutional design to amendment procedures. Third, their exchange paradigm treats politics as mutual exchange of consent rather than top-down command. Government achieves legitimacy when citizens “purchase” its services consensually through constitutional agreement. The consent-holding framework formalizes this metaphor rigorously through stakes-weighted alignment metrics.

Finally, Buchanan and Tullock model optimal decision rules as minimizing total costs combining external costs (harm from decisions affecting you without your consent) and decision costs (time and effort required to reach agreement). Our friction metric  $F(d)$  captures external costs precisely as stakes-weighted deviations from stakeholder ideal points. The framework extends Buchanan and Tullock in four crucial respects. First, we introduce stakes-weighting  $s_i(d)$ , recognizing that individuals are heterogeneously affected by policies. Second, while Buchanan focuses on one-time constitutional founding moments, we model consent-holding as continuously operating through  $H_t(d)$ , tracking legitimacy dynamically as institutional configurations evolve. Third, Buchanan discusses “the” social contract; we specify that consent-holding varies across domains  $d$ , with different optimal structures for taxation, criminal justice, environmental regulation, and community norms. Fourth, Buchanan provides normative theory; we operationalize concepts through  $\alpha(d, t)$  and  $F(d, t)$ , enabling empirical validation of constitutional designs rather than purely philosophical justification.

## 2.2 Social Choice Theory and Impossibility Results

Building on [Dahl \(1956\)](#) foundational analysis of democratic theory showing how pluralist democracy requires balancing majority rule with minority rights, Arrow’s [\(1951\)](#) impossibility theorem establishes that no ranked voting system can simultaneously satisfy four seemingly minimal desiderata: Pareto efficiency, non-dictatorship, independence of irrelevant alternatives, and unrestricted domain. This result demonstrates that perfect democratic aggregation is mathematically impossible, not merely practically difficult. The Gibbard-Satterthwaite theorem ([Gibbard, 1973](#); [Satterthwaite, 1975](#)) extends this impossibility to strategy-proofness: any non-dictatorial voting mechanism over three or more alternatives is manipulable.

Recent quantitative versions show these aren’t merely theoretical concerns but quantifiably common, with [Keller \(2012\)](#); [Mossel et al. \(2012\)](#); [Friedgut et al. \(2011\)](#) quantifying how often Arrow’s impossibility manifests in real voting scenarios with finite electorates. Sen’s [\(2017\)](#) expanded treatment of collective choice integrates economics and ethics, introducing the capability approach that maps di-

rectly onto our effective voice concept. Sen (1999) argues in *Development as Freedom* that development should be measured not by utility or resources alone but by capabilities—freedoms to achieve valued functionings like health, education, and political participation. This provides theoretical grounding for our  $\text{eff\_voice}_i(d)$  term: possessing formal consent power  $C_i > 0$  without resources, education, or political freedom represents low capability.

The consent-holding framework relates to social choice theory as meta-analysis rather than competitor. Where Arrow and Gibbard-Satterthwaite ask “which aggregation rule is best?”, we ask “how legitimate is any given aggregation rule?” This shift has three implications. First, our framework doesn’t compete with impossibility results; it builds on them by providing tools for measuring consequences of unavoidable trade-offs. Since perfect rules don’t exist, we need metrics for comparing imperfect options. Second, stakes-weighting  $s_i(d)$  isn’t present in classical social choice theory, which typically assumes equal weights. This extension allows domain-specific analysis: simple majority may be optimal for low-stakes routine legislation, while supermajority or even consensus becomes appropriate when stakes concentrate heavily.

Our stakes-weighting approach builds on but diverges from weighted voting power analysis (Banzhaf III, 1965; Shapley and Shubik, 1954). The Shapley-Shubik and Banzhaf power indices measure *effective* voting power given formal weights in committee systems—recognizing that a voter with 40% weight may have more than 40% actual power if they’re pivotal in coalitions. This literature addresses measurement of power distribution within given institutional arrangements. Our framework addresses the prior question: how should consent power  $C_i, d$  be allocated based on stakes  $s_i(d)$ ? While power index theory takes weights as given and calculates resulting influence, we propose stakes as foundation for determining appropriate weights. Future work integrating these approaches could specify stakes-weighted allocations, then apply Banzhaf or Shapley-Shubik indices to measure resulting effective voice, combining normative allocation principles with positive power analysis.

## 2.3 Stakeholder Theory and Corporate Governance

Building on Pitkin (1967) foundational work on representation distinguishing substantive versus descriptive representation and acting for constituents, Freeman’s (1984) stakeholder approach argues that firms should create value for all stakeholders—employees, suppliers, communities, customers, shareholders—not just maximize shareholder returns. This challenges Friedman (1970) shareholder primacy doctrine, which treats profit maximization as the sole corporate responsibility and argues that corporate social responsibility beyond shareholder wealth maximization is fundamentally misguided. The 2019 Business Roundtable statement endorsing stakeholder capitalism, signed by 200 CEOs, marks mainstream acceptance of Freeman’s stakeholder view (Business Roundtable, 2019), representing a significant shift from the Friedman doctrine. Phillips (2003) further develops this framework by distinguishing stakeholders by the moral obligation owed to them versus their ability to affect the organization, providing a typology of stakeholder legitimacy that maps onto our stakes-consent framework.

The framework operationalizes Freeman’s insights by defining stakeholders precisely as agents with  $s_i(d) > 0$  in corporate domains. Current governance structures grant consent power almost exclusively to shareholders: they elect boards, approve major transactions, and receive residual claims. Employees, despite high stakes in employment security, working conditions, and workplace norms, hold negligible  $C_i$  in most Anglo-American firms. This generates low  $\alpha(d_{\text{corporate}})$  when stakes are calculated comprehensively. The framework predicts such misalignment produces friction  $F(d)$ : labor disputes, regulatory pressures, reputation damage, difficulty attracting talent.

Comparative corporate governance research validates these predictions. Vitols (2011) documents

how German codetermination—mandatory worker representation on supervisory boards—constrains hostile takeovers and maintains stakeholder orientation. Workers’ voice (high  $\alpha_{workers}(d)$ ) prevents zero-sum shareholder maximization strategies. Fauver and Fuerst (2011) show codetermined firms invest more in worker training and career development; higher  $\alpha$  produces performance improvements in human capital domains.

## 2.4 Common-Pool Resource Governance

Ostrom’s (1990) groundbreaking work on common-pool resources challenges both “tragedy of the commons” pessimism and top-down state solutions. Through field studies of fisheries, forests, irrigation systems, and groundwater basins across continents, she demonstrates that resource users frequently develop effective self-governance without privatization or centralized authority. Her eight design principles for successful commons management include particularly relevant insights for consent-holding theory. Design Principle 3 requires that “most individuals affected by the operational rules can participate in modifying the operational rules”—essentially mandating high  $\alpha(d_{rules})$  for those with high  $s_i(d_{resources})$ . Design Principle 8 specifies nested enterprises for larger systems, enabling polycentric governance with consent-holding at multiple scales. This builds on earlier insights from Ostrom et al. (1961) on polycentric systems, demonstrating how multiple governing authorities at different scales can achieve better outcomes than monocentric alternatives.

The framework formalizes Ostrom’s intuitions. Her “collective choice arrangements” represent  $H_t(d)$  mappings where users participate in rule modification. Her design principles can be reinterpreted as conditions enabling high  $\alpha(d)$ : clear boundaries (defining who holds  $s_i$ ), local monitoring (ensuring  $C_i$  holders possess information), graduated sanctions (responses to low- $\alpha$  violations), and conflict resolution mechanisms (managing  $F(d)$  when it arises). Successful commons maintain high consent alignment; failed commons exhibit persistent misalignment between stakes and voice.

Recent empirical work validates this interpretation quantitatively. Cox et al. (2010) conduct a meta-analysis showing that Ostrom’s design principles predict commons sustainability across diverse contexts, providing systematic evidence that consent alignment mechanisms enable effective resource governance. Yadav et al. (2021) analyze 83 Amazonian communities managing arapaima fisheries, showing that Ostrom’s design principles predict ecological outcomes systematically. Communities exhibiting collective choice arrangements (high  $\alpha$ ) maintain sustainable fish stocks; those lacking such arrangements experience depletion.

## 2.5 Deliberative Democracy and Mini-Publics

Building on Dahl (1971) polyarchy framework of participation and opposition and Mill (1861) considerations on representative government balancing participation and competence, Habermas’s (1984; 1990) communicative action theory distinguishes strategic action (oriented toward achieving one’s goals) from communicative action (oriented toward mutual understanding through reasoned argument). Legitimate norms are those acceptable to all affected parties through rational discourse free from coercion. His discourse principle holds that “only those norms can claim validity that could meet with the acceptance of all concerned in their capacity as participants in a practical discourse.” This maps onto consent-holding directly: “all concerned” represents our affected set  $S_d = \{i | s_i(d) > 0\}$ , while “acceptance” requires  $C_i > 0$  in decision procedures  $H_t(d)$ .

Fishkin’s (2009; 2018) deliberative polling research operationalizes these theoretical commitments. By convening randomly selected representative samples, providing balanced information, facilitating structured deliberation, and measuring preference changes, deliberative polls demonstrate that informed

public judgment shifts significantly through discourse. Citizens’ assemblies extend deliberative innovation to consequential policy domains. The Irish Citizens’ Assembly (2016-2018) addressed abortion and climate change through 99 randomly selected citizens deliberating after expert input, demonstrating how sortition combined with deliberation can shift preferences systematically (Farrell et al., 2019). Courant and Bourgeron (2021) analyze the French Citizens’ Convention on Climate (2019-2020), which generated 149 policy proposals from 150 randomly selected participants through sortition and deliberation, with many subsequently adopted into legislation.

The framework interprets these innovations as institutional experiments raising  $\alpha(d)$  through sortition and deliberation. Random selection approximates equal  $C_i$  for participants; demographic stratification can approximate stakes-weighting if groups correlate with  $s_i(d)$ . Learning phases improve  $\text{eff\_voice}_i$  through information provision; deliberation structures enable preference refinement.

Our stakes-weighted consent framework confronts democratic equality arguments directly. Building on Mill (1859) foundations regarding individual liberty, consent, and limits of state power, Christiano (2008) defends equal political voice on dignity grounds: each person possesses equal moral status, entitling them to equal say in collective decisions regardless of stakes or competence. Waldron (1999) argues that persistent disagreement about what justice requires makes equal voice procedurally fair even if some possess superior judgment. Brighouse and Fleurbaey (2010) examine whether proportional influence could improve democratic outcomes but conclude that equal voice better respects equality of persons.

We acknowledge this tension while distinguishing *political* domains from *governance* domains generally. In constitutional fundamentals and citizenship rights, equal voice may be intrinsically required by equal moral status—each person gets one vote precisely because they are persons, not because they possess equal stakes. But many governance domains are not *political* in this sense: corporate boards allocating firm resources, technical committees setting safety standards, platform algorithms moderating speech, common-pool resource users managing fisheries. In these contexts, stakes-weighting may be both more efficient (reducing friction, improving outcomes) and more legitimate (those bearing consequences should influence decisions proportionally). The framework enables empirical testing: do equal-voice or stakes-weighted mechanisms generate higher measured legitimacy  $L(d, t)$  in different domain types?

## 2.6 Algorithmic Governance and Platform Legitimacy

Grimmelikhuijsen et al. (2022) identify three legitimacy dimensions for algorithmic decision-making: input (did citizen preferences inform design?), throughput (does the algorithm follow fair procedures?), and output (do outcomes align with public values?). Current algorithmic governance exhibits severe deficits across all three dimensions. Citizens rarely participate in algorithm design (low input legitimacy), decision-making processes remain opaque black boxes (low throughput legitimacy), and outcomes often replicate historical discrimination (questionable output legitimacy). Kleinberg et al. (2017) demonstrate inherent trade-offs in fair determination of risk scores, showing that multiple incompatible definitions of algorithmic fairness exist—making it impossible to satisfy all fairness criteria simultaneously, analogous to Arrow’s impossibility theorem in social choice.

Waldman and Johnson (2022) show that high-stakes algorithmic decisions (healthcare allocation, criminal sentencing) are perceived as less legitimate than human decisions even when outcomes are identical. The consent-holding framework diagnoses these challenges structurally. Algorithmic decision-making creates domains  $d_{\text{algorithm}}$  where algorithms or their designers hold  $C \approx 1$  while affected citizens have  $C \approx 0$  despite high  $s_i(d)$ . Credit scoring algorithms determine loan access (high  $s_i$  for applicants);

hiring algorithms control employment opportunities (high  $s_i$  for candidates); content moderation algorithms shape speech norms (high  $s_i$  for platform users). In each case, current  $\alpha(d_{algorithm}) \approx 0$  because high-stakes populations are excluded from  $H_t(d)$ .

Platform responses attempting to raise  $\alpha$  reveal understanding of legitimacy deficits. Meta’s Oversight Board provides independent content moderation appeals, slightly raising  $\alpha(d_{moderation})$  by giving users contestation rights, though Douek (2022) notes this provides only limited voice expansion while maintaining corporate control over fundamental rules. YouTube Creator Councils consult high-profile creators, extending partial  $C_i$  to stakeholders whose  $s_i$  is highest.

## 2.7 Voting Power Indices and Coalition Analysis

The power indices literature demonstrates that voting weight  $\neq$  voting power. Banzhaf III (1965) measures critical voter frequency: how often removing your vote changes outcomes from win to loss. Shapley and Shubik (1954) measure pivotal voter frequency in sequential coalition formation. Felsenthal and Machover (1998) provide comprehensive comparison of these approaches, demonstrating that Banzhaf and Shapley-Shubik indices often diverge substantially and measure different aspects of voting power. These indices often diverge dramatically from nominal weights—Germany holds the most European Council votes but doesn’t possess proportional power due to coalition dynamics. Similar phenomena arise in corporate boards (blockholders vs. minority shareholders), legislatures (swing voters vs. party leaders), and qualified majority systems (Security Council veto players).

These insights directly inform consent-holding operationalization. Naive approaches measure  $C_i$  as voting weight (shares held, seats controlled). Sophisticated approaches use power indices accounting for coalition structures. In weighted voting contexts (shareholders, federalism), qualified majority rules (constitutional amendments), and veto player systems (UN Security Council), indices capture actual influence more accurately than nominal weights.

The framework integrates power indices into legitimacy measurement:  $\alpha(d, t) = \frac{\sum_i s_i(d) \cdot \text{PowerIndex}_i(d, t)}{\sum_i s_i(d)}$ , where  $\text{PowerIndex}_i$  represents Banzhaf, Shapley-Shubik, or domain-appropriate measures. This refinement matters most when vote concentration enables blocking coalitions. Consider corporate governance: a minority shareholder with 20% equity plus veto rights over major transactions wields power far exceeding their ownership share. Measuring  $C_i = 0.20$  understates influence; calculating Banzhaf index accounting for veto power provides accurate assessment.

Recent extensions analyze endogenous coalition formation (Aumann and Myerson, 1988), showing how equilibrium structures emerge from bargaining. This connects to consent-holding’s dynamic aspect:  $H_t(d)$  evolves as agents form alliances, shifting power distributions. Nash bargaining solutions (Nash, 1950) maximize products of utility gains subject to Pareto efficiency—structurally similar to stakes-weighted consent maximization. Kalai and Smorodinsky (1975) propose alternative axiomatizations highlighting trade-offs between equality (proportional gain-sharing) and efficiency (Pareto optimality), demonstrating solution multiplicity absent unique normative commitments—precisely what Theorem 3 predicts.

## 2.8 Relational Autonomy and Consent Capacity

Mackenzie (2014) three-dimensional autonomy framework distinguishes self-determination (choosing one’s own life path), self-governance (regulating one’s actions), and self-authorization (taking responsibility for choices). Traditional liberal autonomy assumes atomistic individuals; relational approaches (Mackenzie and Stoljar, 2000) recognize that autonomy is socially constituted—relationships and social structures fundamentally enable or constrain autonomous choice rather than merely influencing



pre-existing capacities. Nedelsky (1989) analyzes how oppressive social structures systematically constrain women’s autonomy through relational mechanisms, demonstrating that coercion operates not only through direct force but through systematic limitation of available choices. Oppressive systems constrain capacity for self-governance—gender oppression limits women’s educational access, economic opportunities, and freedom from violence, directly undermining autonomous choice.

Koggel (2022) extends this to global justice, arguing that respecting autonomy requires enabling threshold capabilities, not merely non-interference. Autonomy necessitates freedom conditions: political liberties (speech, association, conscience) and personal liberties (movement, bodily autonomy, freedom from violence). Agents lacking these conditions cannot exercise meaningful consent even if formally included in  $H_i(d)$ .

These insights address the framework’s handling of  $\text{eff\_voice}_i$ . Relational autonomy equals effective voice in our legitimacy equation. Simply granting  $C_i > 0$  (voting rights) without resources, education, or freedom produces low  $\text{eff\_voice}_i$ —formal authority without capacity to exercise it. Oppressive structures systematically reduce both stakes recognition (dominant groups deny subordinated groups’  $s_i$ ) and consent power (exclusion from  $H_i(d)$  even for high-stakes domains).

This perspective addresses three framework challenges. First, it resolves circularity concerns: “Who decides who’s in  $H_i(d)$ ?” Answer: those with stakes plus capacity, considering relational constraints that may undermine apparent consent. Second, it handles vulnerable populations ethically. Proxy consent becomes necessary when capacity is impaired, but structures should enable gradual inclusion as capability develops rather than permanent exclusion. Third, it enables justice analysis: systematic exclusion of groups with high  $s_i$  but low  $\text{eff\_voice}_i$  constitutes legitimacy deficit diagnosable through  $\alpha(d)$  measurement.

Application to research ethics illustrates these dynamics. Standard approaches grant legal guardians consent authority over cognitively impaired individuals. Relational approaches recognize impaired persons retain partial capacity and value particular relationships beyond legal guardianship—an older sibling may understand needs better than distant legal guardians. The framework prescription: allocate partial  $C_i$  based on measured capacity and expand  $H_i(d)$  to include chosen trusted relationships, raising  $\alpha(d_{\text{research}})$  for the affected individual.

### 3 Formal Framework: Primitives, Axioms, and Theorems

This section establishes the framework’s formal foundations through precise definitions, minimal axioms, and structural theorems. The approach proceeds deductively: from spare assumptions about collective decision-making to necessary conclusions about consent-holding’s existence, friction’s inevitability, and legitimacy’s measurement.

#### 3.1 Primitives and Definitions

We begin with foundational concepts requiring no prior theoretical commitment. An **agent** is any entity capable of selecting among actions, indexed  $i \in A = \{1, \dots, N\}$ . Agents may be individuals, organizations, algorithms, or collective bodies—the framework remains agnostic about internal composition. A **domain** represents a decision-relevant sphere—a policy area, firm process, household choice, or any context requiring action selection. The set of domains is  $D = \{d_1, \dots, d_M\}$ . Each domain  $d$  admits a set of possible actions  $X_d$ , from which one action  $x_d \in X_d$  must be selected.

**Outcomes** represent realized states resulting from action vectors  $\mathbf{x} = (x_{d_1}, \dots, x_{d_M})$  through an environment mapping  $E : \prod_d X_d \rightarrow O$ , where  $O$  denotes the outcome space. An agent  $i$ ’s **stake** in domain  $d$ , denoted  $s_i(d) \geq 0$ , quantifies sensitivity to outcomes in that domain. Stakes may reflect material

exposure, legal consequences, capability impacts, or existential threats.

Each agent possesses **preferences** over outcomes, represented either as complete orderings  $\succeq_i$  or utility functions  $U_i : O \rightarrow \mathbb{R}$ . Preferences induce ideal points  $x_{i,d}^*$  in each domain—the action agent  $i$  most prefers given others’ anticipated choices.

**Consent** represents the normative right to decide in a domain—who may authoritatively say “yes” or “no” to proposed actions. Following **Locke (1980)** consent theory foundations, political obligation derives from voluntary agreement; actual consent is required for legitimate authority. The **consent-holder mapping**  $H_t(d) \in \Delta(C)$  specifies the distribution of decision authority over possible holders  $C$  at time  $t$ . Individual **consent power**  $C_{i,d} \in [0, 1]$  represents agent  $i$ ’s effective share of decision authority in domain  $d$ , with  $\sum_i C_{i,d} = 1$ .

### 3.2 Axioms

The framework rests on seven axioms representing minimal commitments about collective decision-making.

**A1. Action Precedence:** Every non-null outcome in a domain is produced by some action (including “do nothing”).

**A2. Decision Requirement:** Every action is selected by some decision procedure (choice, rule, randomization, delegation).

**A3. Shared Reality:** Outcomes alter a world co-occupied by multiple agents; externalities exist.

**A4. Finitude:** Agents have finite time, attention, and cognitive capacity; no single agent can decide everything alone.

**A5. Plurality:** Agents’ preference orderings differ on at least some domains.

**A6. Salience:** For each domain, at least one agent has  $s_i(d) > 0$ .

**A7. Fallibility/Subjectivity:** Perception and valuation are frame-dependent; no universal content-level value ordering is logically forced.

### 3.3 Theorem 1: Consent-Holding Necessity

**Theorem 3.1** (Consent-Holding Necessity). *In any domain  $d$  where a non-null outcome occurs, there exists a consent-holder mapping  $H_t(d)$ .*

*Proof Sketch.* By A1-A2, any outcome resulted from an action selected through some procedure. A procedure implies a locus of control—the entity/entities choosing the action, establishing the choice rule, or delegating to randomization. This locus constitutes  $H_t(d)$ . Therefore, denying  $H_t(d)$ ’s existence contradicts A2. ■ ■

### 3.4 Theorem 2: Inevitable Friction

**Theorem 3.2** (Inevitable Friction). *If there exist agents  $i, j$  with divergent preferences on domain  $d$  and  $s_i(d), s_j(d) > 0$ , then unless  $H_t(d)$  exactly reproduces stakes-weighted unanimity, at least one agent experiences moral/political friction.*

We formalize **friction** in domain  $d$  as:

$$F(d, t) = \sum_i s_i(d) \cdot \delta(x_d(t), x_{i,d}^*) \quad (1)$$

where  $x_{i,d}^*$  represents agent  $i$ ’s ideal action and  $\delta$  measures divergence. For discrete choices,  $\delta(x, x^*) = 0$  if  $x = x^*$ , else 1. For continuous policy spaces,  $\delta(x, x^*) = |x - x^*|$  captures distance from ideal points.

Introducing tolerance thresholds  $\tau_i$  yields:

$$F_\tau(d, t) = \sum_i s_i(d) \cdot \max(0, \delta(x_d, x_{i,d}^*) - \tau_i) \quad (2)$$

### 3.5 Definition 1: Legitimacy as Consent Alignment

We operationalize legitimacy through stakes-weighted consent alignment. Define the **affected set**  $S_d = \{i | s_i(d) > 0\}$ . **Consent alignment** is:

$$\alpha(d, t) = \frac{\sum_{i \in S_d} s_i(d) \cdot \text{eff\_voice}_i(d, t)}{\sum_{i \in S_d} s_i(d)} \quad (3)$$

where  $\text{eff\_voice}_i$  represents agent  $i$ 's effective decision power in  $H_t(d)$ .

**Definitional Note:** This is a *measurement framework*, not a derived result. We define legitimacy as the degree to which consent power tracks stakes distribution, making the concept empirically tractable. The framework's predictive power lies in the hypothesis that low  $\alpha$  generates observable friction—a claim requiring empirical validation beyond the definition itself.

A minimal procedural legitimacy condition requires  $\alpha(d, t) \geq \tau$  for society-specific threshold  $\tau$ . Persistent  $\alpha < \tau$  predicts observable friction through unrest, exit, sabotage, or normative decay.

### 3.6 Postulate 1: Competence-Consent Trade-Off

We model overall legitimacy as combining consent alignment and performance:

$$L(d, t) = w_1 \cdot \alpha(d, t) + w_2 \cdot P(d, t) \quad (4)$$

where  $\alpha(d, t)$  represents stakes-weighted consent alignment,  $P(d, t)$  denotes performance/competence metrics, and  $w_1, w_2 \geq 0$  reflect society-specific weights on voice versus results.

This is a **postulated relationship** rather than a derived theorem. The linear combination assumes legitimacy trades off between consent and competence, but alternative functional forms (multiplicative, threshold-based) are possible. Empirical work validating this specification against alternatives remains a key research agenda.

**Remark on Weight Determination:** The weights  $w_1, w_2$  are not free parameters requiring external normative specification, but endogenous scope conditions revealed through constitutional-level decisions (see Section 10.1 for full meta-legitimacy resolution). Societies whose weight configurations produce excessive friction face structural pressure to reform. We can characterize admissible weight functions axiomatically: any stable society must satisfy  $w_2/w_1 > f(\text{Var}[s_i(d)])$  where  $f$  is a function of stakes heterogeneity derived from friction minimization. Future empirical work will estimate weights via:

$$(w_1^*, w_2^*) = \arg \min_{w_1, w_2} \mathbb{E}[F(d, t; w_1, w_2)] \quad (5)$$

where friction minimization across constitutional reforms provides revealed preference data for weight estimation. Sequential Monte Carlo methods (Lux and Schiffko, 2018) enable parameter estimation for agent-based models through particle filtering, providing techniques applicable to calibrating consent-holding frameworks against empirical institutional data. This transforms weight determination from a normative choice into an empirical optimization problem, sidestepping the meta-legitimacy regress.

This formulation makes explicit that different systems optimize different points on the legitimacy frontier. Technocracies maximize  $P$ , often sacrificing  $\alpha$  by concentrating consent in experts. Direct



democracies maximize  $\alpha$  through universal suffrage, potentially reducing  $P$  on technical domains where distributed knowledge is sparse.

### 3.7 Theorem 3: Minimal Absolutism from Relativism

**Theorem 3.3** (Relativism  $\Rightarrow$  Minimal Absolutism). *Given A7 (value frame-dependence), the claim “all value judgments are frame-relative” is coherent only if the **structure** enabling frames is invariant. Therefore, at least one absolute exists: the necessity of consent-holding over shared outcomes wherever A1-A6 hold.*

*Proof Sketch.* Suppose all value claims are frame-dependent (A7). Frame-dependence presupposes frames exist—perspectives from which valuations occur. Frames belong to agents inhabiting shared reality (A3) with plural preferences (A5). These agents make decisions affecting each other (A1-A2). Such decisions require consent-holder mappings  $H_i(d)$  (Theorem 1). Therefore, relativism about content-level values doesn’t extend to structural necessities. ■ ■

## 4 Operationalization: Empirical Measurement and Identification

The theoretical framework provides analytical tools for understanding consent-holding structures. This section bridges theory and empirical application by specifying how the framework’s core concepts can be measured, how causal relationships can be identified econometrically, and what testable predictions emerge.

### 4.1 Formal Measurement Framework

We operationalize consent-holding through a consent matrix  $\mathbf{C} \in [0, 1]^{N \times M}$ , where each element  $C_{i,d}$  represents agent  $i$ ’s effective decision share in domain  $d$ , subject to the normalization constraint  $\sum_i C_{i,d} = 1$ . This matrix captures both de jure authority and de facto power. In simple majority voting systems with equal suffrage,  $C_{i,d} = 1/N_{voters}$  for all enfranchised  $i$  and  $C_{i,d} = 0$  for excluded populations. In shareholder governance,  $C_{i,d} = \text{shares}_i / \text{total\_shares}$ . In technocratic systems,  $C_{i,d} = 1/|E|$  if agent  $i$  belongs to the expert set  $E$ , zero otherwise.

Complementing the consent matrix, the stakes vector  $\mathbf{s}(d) \in \mathbb{R}_{\geq 0}^N$  quantifies each agent’s exposure to consequences in domain  $d$ . Stakes measurement presents both conceptual and practical challenges. Conceptually, stakes may reflect material exposure (tax burden relative to income), capability impacts (health outcomes affected), or existential threats (survival risks from climate policy). Different domains may legitimately employ different stakes conceptions.

Combining these elements, consent alignment in domain  $d$  at time  $t$  is measured as:

$$\alpha(d, t) = \frac{\sum_{i \in S_d} s_i(d) \cdot \text{eff\_voice}_i(d, t)}{\sum_{i \in S_d} s_i(d)} \quad (6)$$

where  $S_d = \{i | s_i(d) > 0\}$  denotes the affected set and  $\text{eff\_voice}_i$  represents agent  $i$ ’s effective decision power accounting for both formal authority  $C_{i,d}$  and capacity constraints.

### 4.2 Friction Metrics and Tolerance-Weighted Extensions

Political friction represents the stakes-weighted aggregate deviation between realized outcomes and stakeholder preferences. In its basic form:

$$F(d, t) = \sum_i s_i(d) \cdot \delta(x_d(t), x_{i,d}^*) \quad (7)$$

For continuous policy spaces, Euclidean distance  $\delta(x, x^*) = |x - x^*|$  captures proximity to ideal points. The tolerance-weighted friction measure incorporating agent-specific tolerance parameters  $\tau_i \geq 0$  is:

$$F_\tau(d, t) = \sum_i s_i(d) \cdot \max(0, \delta(x_d(t), x_{i,d}^*) - \tau_i) \quad (8)$$

This captures that agents tolerate “good enough” governance within zones of acceptability, mobilizing only when deviations exceed tolerance thresholds.

### 4.3 Empirical Identification Strategies

The core empirical prediction connecting alignment to friction generates testable hypotheses through panel regression specifications:

$$F_{d,t} = \beta_0 + \beta_1 \cdot \alpha_{d,t} + \beta_2 \cdot P_{d,t} + \gamma \cdot X_{d,t} + \mu_d + \lambda_t + \varepsilon_{d,t} \quad (9)$$

where  $F_{d,t}$  represents friction,  $\alpha_{d,t}$  denotes consent alignment,  $P_{d,t}$  captures performance outcomes,  $X_{d,t}$  includes control variables,  $\mu_d$  represents domain fixed effects,  $\lambda_t$  represents time fixed effects, and  $\varepsilon_{d,t}$  is the error term.

The framework’s theoretical predictions constrain coefficient signs:  $\beta_1 < 0$  (higher alignment reduces friction),  $\beta_2 < 0$  (better performance reduces friction). Instrumental variable strategies address endogeneity concerns by exploiting exogenous variation in consent structures. Historical franchise expansions driven by international diffusion provide quasi-experimental variation.

### 4.4 Testable Predictions and Empirical Hypotheses

**Hypothesis 1 (Alignment-Friction Relationship):** Across domains and time periods, higher consent alignment  $\alpha(d, t)$  predicts lower friction  $F(d, t+k)$  with lags  $k$  reflecting institutional adjustment speeds:

$$\frac{\partial F(d, t+k)}{\partial \alpha(d, t)} < 0 \quad \text{for } k \geq 0$$

**Hypothesis 2 (Stakes-Consent Covariance):** Institutional reforms increasing the covariance between stakes and consent power— $\text{Cov}(s_i(d), C_{i,d})$ —reduce friction through alignment improvement.

**Hypothesis 3 (Threshold Effects):** Domains with alignment below societal tolerance thresholds— $\alpha(d) < \tau_{\text{legitimacy}}$ —exhibit discontinuously higher instability, generating nonlinearity in the alignment-friction relationship.

**Hypothesis 4 (Temporal Dynamics):** Persistent friction  $F(d, t)$  predicts future alignment increases  $\alpha(d, t+k)$  through institutional reform pressure:

$$\frac{\partial \alpha(d, t+1)}{\partial F(d, t)} > 0$$

**Hypothesis 5 (Performance Interactions):** The alignment-friction relationship weakens in domains with high performance  $P(d, t)$ , as competent governance partially compensates for voice deficits.

## 5 Social Contract Theories as Distribution Mechanisms

Social contract theories can be reinterpreted through the consent-holding framework as different proposals for how to allocate consent power  $C_i$  across agents in various domains. Rather than treating these theories as competing comprehensive doctrines, we analyze them as institutional design proposals optimizing different legitimacy functions subject to domain-specific constraints.

### 5.1 Rawlsian Justice as Maximin Consent

Rawls’s (1971) difference principle can be formalized as maximizing the minimum effective voice:

$$\max_{C_{i,d}} \min_i \{\text{eff\_voice}_i(d)\} \quad (10)$$

subject to basic liberties constraints ensuring  $C_{i,d} > 0$  for all citizens in political domains. This generates predictions about institutional design: political equality (one person, one vote) in constitutional domains, economic redistribution raising least-advantaged citizens’ capability to exercise voice, and priority rules protecting basic liberties even when aggregate welfare would benefit from violation.

### 5.2 Utilitarian Consent as Weighted Aggregation

Classical utilitarianism maximizes stakes-weighted welfare:

$$\max_{x_d} \sum_i s_i(d) \cdot U_i(x_d) \quad (11)$$

This doesn’t directly specify consent allocation, but combined with epistemic assumptions that affected parties possess superior information about their own stakes  $s_i(d)$ , it motivates giving consent power proportional to stakes—exactly our  $\alpha(d)$  alignment measure. The framework reveals utilitarianism’s implicit consent structure: let those with stakes decide, weighted by their exposure.

### 5.3 Libertarian Consent as Property Rights

Nozickean libertarianism allocates consent power through property rights:  $C_{i,d} = 1$  if domain  $d$  involves only  $i$ ’s property, distributed according to ownership shares otherwise. This generates high  $\alpha(d)$  for domains where property rights align with stakes (personal consumption choices) but potentially low  $\alpha$  for domains with externalities (pollution, network effects) where those holding property rights differ from those bearing consequences.

The framework doesn’t adjudicate between these theories normatively but provides tools for comparing their institutional predictions and empirical performance across domains.

## 6 Historical Validation: Case Studies in Consent Alignment Dynamics

The framework’s predictive power rests on historical validation. We examine seven domains where consent alignment  $\alpha(d)$  evolved over time, generating observable friction  $F(d)$  when misaligned and stability when aligned. Each case demonstrates the framework’s core prediction: persistent low  $\alpha(d)$  generates escalating friction until institutional reform raises alignment above threshold  $\tau$ , or suppression temporarily contains mobilization.

### 6.1 Suffrage Expansion: Gradual Consent Broadening

Women’s suffrage movements (1890s-1970s) demonstrate the predicted  $\alpha$ - $F$  dynamics. Building on foundational principles articulated in Stanton et al. (1848) Declaration of Sentiments claiming equal political rights, women held extreme stakes in political domains affecting family law, property rights, employment regulation, and reproductive policy ( $s_{\text{women}}(d) \gg 0$ ), yet possessed zero formal consent power ( $C_{\text{women}} = 0$ ) until franchise extension. This generated high friction: suffragist organizing, civil disobedience, protest movements. New Zealand (1893), Australia (1902), Finland (1906), and Norway (1913) extended franchise early; the United States (1920), United Kingdom (1928), France (1944), and Switzerland (1971) delayed decades longer. Teele (2018) demonstrates that electoral logic drove gradual

enfranchisement in the United States, with competitive mobilization among political parties accelerating expansion as friction intensified.

The framework predicts that earlier adopters experienced lower friction costs from exclusion—smaller suffrage movements, less civil unrest. Later adopters faced escalating friction as international demonstration effects raised women’s consciousness of exclusion. [Ramirez et al. \(1997\)](#) document how suffrage movements formed transnational networks, with international diffusion accelerating adoption as demonstration effects intensified across national boundaries. Empirical validation could test whether protest intensity  $F(d, t)$  correlates negatively with time-to-adoption, controlling for other democratization factors.

## 6.2 Abolition Movements: Maximum Stakes, Zero Consent

Enslaved populations held maximal stakes in slavery policy domains ( $s_{\text{enslaved}}(d) = \text{existential}$ )—literally life, liberty, and bodily autonomy—yet possessed zero consent power by definition ( $C_{\text{enslaved}} = 0$ ). This generated extreme misalignment  $\alpha(d_{\text{slavery}}) \approx 0$  despite involving the highest possible stakes.

The framework predicts unsustainable friction: slave rebellions (Haiti 1791-1804, Nat Turner 1831, countless smaller uprisings), abolitionist movements channeling moral friction from sympathetic observers, and ultimately civil war when peaceful adjustment failed (US 1861-1865). [Blackburn \(1988\)](#) documents how the Haitian Revolution and other slave uprisings forced fundamental reconsideration of slavery’s sustainability, demonstrating that high-friction resistance could make exclusionary institutions untenable. Primary sources from enslaved persons like [Equiano \(1789\)](#) provided firsthand documentation of stakes and friction, making the human cost of zero consent visible to broader publics. Britain’s earlier abolition (1833) via compensated emancipation demonstrates an alternative high- $\alpha$  pathway: incorporating enslaved persons’ stakes through proxy representation (abolitionist movements) raised effective  $\alpha$  sufficiently to enable peaceful transition. [Clarkson \(1808\)](#) meticulously documented the abolitionist campaign as a leading researcher, showing how systematic evidence-gathering and mobilization generated friction through moral pressure. Parliamentary advocates like [Wilberforce \(1789–1807\)](#) transformed this grassroots friction into legislative action through decades of speeches and campaigns. [Drescher \(1987\)](#) documents how British abolition succeeded through combining parliamentary lobbying, mass petition campaigns, and sustained moral pressure—effectively raising  $\alpha$  through proxy consent mechanisms before legal emancipation occurred.

## 6.3 Labor Rights and Corporate Codetermination

Workers hold substantial stakes in workplace domains ( $s_{\text{workers}}(d_{\text{workplace}})$  includes employment security, wages, safety, dignity) yet traditionally possessed minimal corporate consent power under shareholder primacy ( $C_{\text{workers}} \approx 0$ ). Early labor organizations like [Knights of Labor \(1878\)](#) articulated stakes claims and demands for worker voice in their foundational principles. This generated labor friction: strikes, unionization drives, regulatory pressure. [Fine \(1969\)](#) documents the 1936-1937 General Motors sit-down strike as the most significant American labor conflict, demonstrating how friction manifestation through direct action forced UAW recognition and fundamentally shifted labor relations.

Different societies responded differently. Germany institutionalized codetermination (1951 Mitbestimmung), granting workers 50% supervisory board representation in large firms—dramatically raising  $\alpha_{\text{workers}}(d_{\text{corporate}})$ . [McGaughey \(2016\)](#) documents how German codetermination emerged from collective bargaining between business and labor during reconstruction periods (1918-1922 and 1945-1951), representing negotiated incorporation rather than revolutionary imposition. The US largely resisted, maintaining low  $\alpha$  through union suppression and shareholder primacy. The framework predicts Ger-

many should exhibit lower ongoing labor friction (fewer strikes, less adversarial labor relations) at cost of potentially lower shareholder returns (lower  $P$  on shareholder-centric metrics). Empirical evidence broadly confirms: Jäger et al. (2022) demonstrate that German firms with codetermination show lower strike rates, longer employee tenure, and stable returns compared to Anglo-American firms, while Fauver and Fuerst (2011); Vitols (2011) document sustained stakeholder orientation and investment in human capital.

#### 6.4 Platform Governance Rebellions (2010s-Present)

Digital platforms create novel consent-holding challenges. Users hold high stakes in content moderation ( $s_i(d_{\text{moderation}})$  includes speech rights, community norms, information access), recommendation algorithms ( $s_i(d_{\text{algorithms}})$  shapes information diet, attention allocation), and governance policy ( $s_i(d_{\text{governance}})$  affects user experience, privacy, monetization). Yet consent power concentrates almost entirely in platform executives and engineers:  $C_{\text{users}} \approx 0$  despite billions affected.

The framework predicts rising friction as stakes grow: #DeleteFacebook movements (2018), advertiser boycotts, regulatory backlash (GDPR 2018, DSA 2022), mass migration to alternatives when they emerge (Twitter/X exodus 2022-2024 to Mastodon, Bluesky, Threads). Gillespie (2018) documents how platform governance legitimacy crises emerge from perceived bias, lack of transparency, and systematic user exclusion from content moderation decisions—classic symptoms of low  $\alpha(d)$  when high-stakes populations lack effective voice. Platforms attempting to raise  $\alpha$  through Oversight Boards (Meta), Creator Councils (YouTube), and community moderation (Reddit) represent elite responses to friction, though these reforms grant only partial voice, maintaining low overall  $\alpha$ .

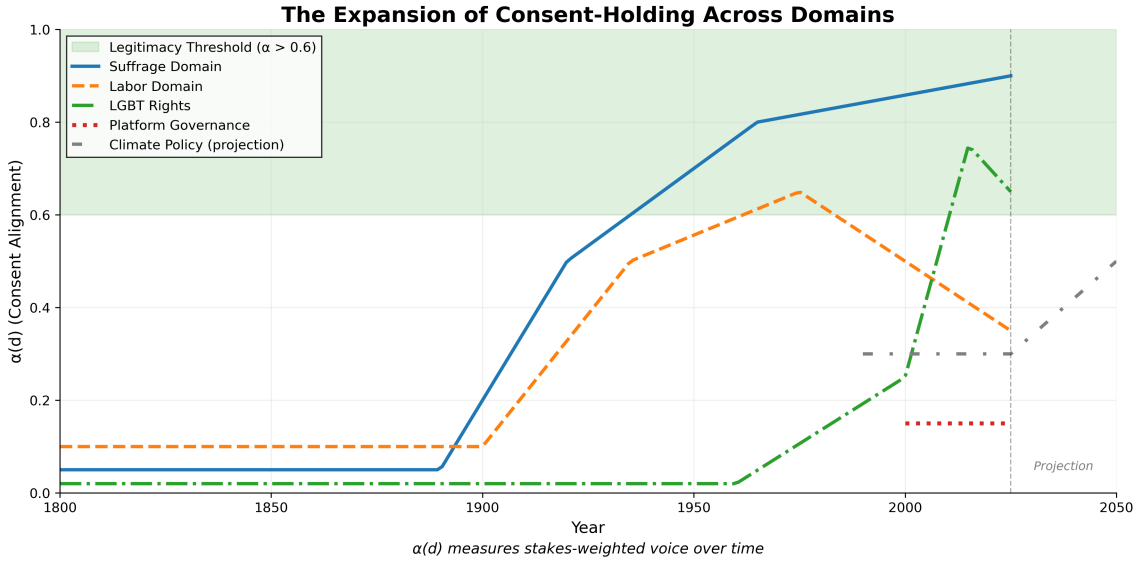


Figure 1: Historical consent alignment trajectories across four domains (suffrage, abolition, labor rights, platform governance) showing predicted dynamics: persistent low  $\alpha$  generates rising friction  $F$  until incorporation or suppression. Suffrage demonstrates gradual incorporation; abolition shows delayed incorporation with violent friction; labor rights exhibits cross-national variation; platform governance shows early-stage friction emergence.

#### 6.5 Scope Conditions: When Friction Fails to Generate Incorporation

Our case selection above focuses on movements that achieved substantial incorporation (suffrage, abolition, labor rights, LGBT rights). However, numerous high-stakes populations have sustained friction

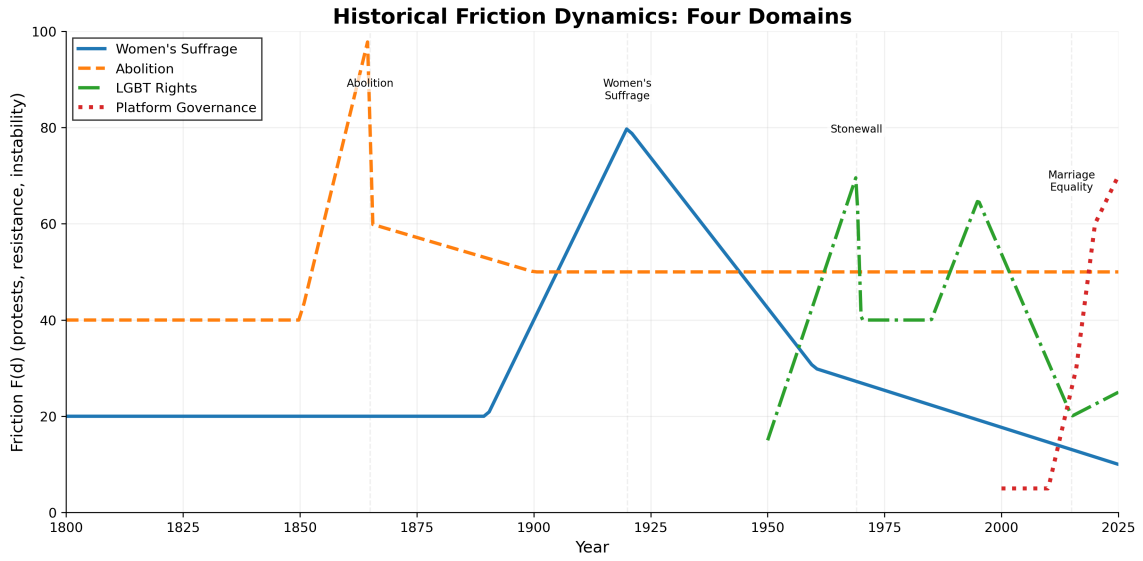


Figure 2: Friction trajectories corresponding to Figure 1. Friction  $F(d, t)$  rises when alignment remains low, spikes during mobilization peaks (suffrage protests 1910s, Civil War 1860s, labor strikes 1930s, platform revolts 2020s), then declines following institutional reform raising  $\alpha$ . Dotted lines indicate counterfactual friction under maintained exclusion.

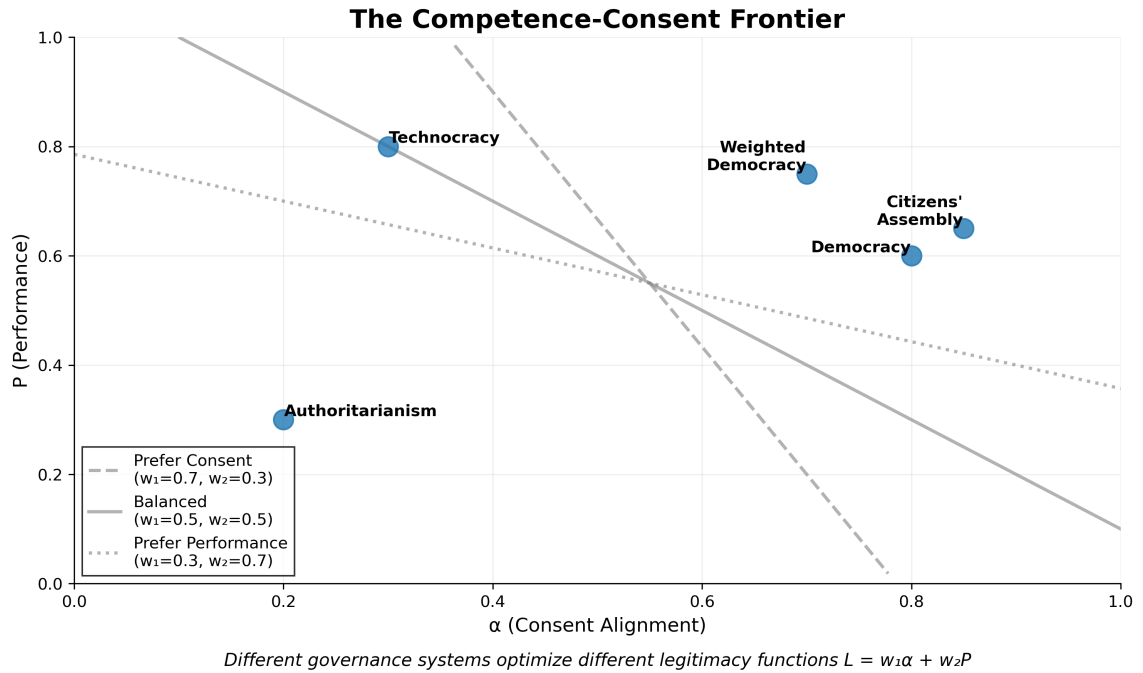


Figure 3: Legitimacy frontier showing trade-offs between consent alignment  $\alpha$  and performance  $P$  across governance systems. Points represent empirical observations: Scandinavian social democracies achieve high  $\alpha$  with moderate-high  $P$ ; technocratic Singapore shows high  $P$  with moderate  $\alpha$ ; failed states exhibit low on both dimensions. The efficient frontier (solid curve) shows maximum achievable  $P$  at each  $\alpha$  level; institutional innovations shift the frontier outward.

for decades or centuries *without* achieving consent expansion, revealing scope conditions requiring theoretical development.

Indigenous sovereignty movements exhibit maximal stakes—land, culture, self-determination—combined with minimal consent power in settler-colonial states, generating sustained friction since colonization. Yet incorporation remains limited despite centuries of mobilization (Lakota/Dakota resistance 1850s-present, Aboriginal land rights struggles in Australia, Māori sovereignty movements in New Zealand). Stateless populations (Rohingya, Palestinians, Kurds) face existential stakes with zero voice in determining their status, producing refugee crises and armed conflict without resolution. Prisoners in most societies have direct stakes in criminal justice policy but negligible voice, despite persistent grievances and occasional rebellions (Attica 1971, UK prison strikes 2016).

These cases suggest friction alone is insufficient for incorporation—elite responses depend on additional factors: (1) *cost of repression* versus accommodation (when repression is cheap relative to consent expansion, suppression persists); (2) *international pressure* and norm diffusion (isolated regimes resist longer than those facing external scrutiny); (3) *coalition availability* among enfranchised groups (reform requires allies with existing voice); (4) *elite interest alignment* with reform (incorporation becomes feasible when elite factions benefit). When repression costs are low, international isolation is possible, and elite interests oppose incorporation, high friction can persist indefinitely through suppression rather than consent expansion.

Future empirical work should model these scope conditions explicitly, predicting when friction generates incorporation versus sustained authoritarianism. This extension would transform the framework from describing  $\alpha$ - $F$  dynamics to predicting trajectories based on initial conditions and contextual parameters.

## 7 Computational Mechanism Comparison: Adaptive Learning Dynamics

Beyond historical validation, we compare consent allocation mechanisms through computational simulation. Monte Carlo experiments (Robert and Casella, 2004) vary mechanisms across diverse preference distributions under adaptive learning dynamics to assess relative performance, using repeated random sampling to obtain numerical results about mechanism performance.

**Methodological Transparency:** The Bayesian learning model implements preference adaptation toward observed outcomes, which by construction reduces friction over time as agents’ ideal points converge toward policy decisions. This is not a test of whether friction *can* reduce—that follows definitionally from preference convergence—but a **comparative test** of which consent allocation mechanism produces superior alignment trajectories when agents adapt to institutional performance. The simulation addresses: given plausible behavioral assumptions (agents learn from outcomes), which mechanism best matches stakeholder interests with institutional decisions?

This computational exercise demonstrates mechanism rankings under specific auxiliary assumptions (Bayesian learning with stakes-weighted attention) rather than validating the framework’s core theoretical claims. The framework’s definitional structure (legitimacy as stakes-weighted alignment, friction as preference deviation) means the Monte Carlo provides illustrative comparison rather than empirical proof. Agent-based modeling (Epstein, 2006) enables bottom-up simulation of complex social systems through individual agent interactions, providing a natural methodology for exploring consent-holding dynamics computationally.



## 7.1 Simulation Design

We simulate 1000 societies, each with  $N = 100$  agents making decisions in  $M = 10$  domains over  $T = 50$  time periods. Agent stakes  $s_i(d)$  are drawn from heterogeneous distributions: some domains exhibit concentrated stakes (e.g., environmental policy affecting coastal residents heavily, inland minimally), others show uniform stakes (e.g., monetary policy affecting all). Preferences  $x_{i,d}^*$  are initially drawn from various distributions (normal, bimodal, skewed) to test robustness.

We compare five consent allocation mechanisms:

1. **Equal voice** (pure democracy):  $C_{i,d} = 1/N$  for all  $i, d$
2. **Stakes-weighted**:  $C_{i,d} = s_i(d) / \sum_j s_j(d)$
3. **Random** (sortition):  $C_{i,d} = 1$  for randomly selected  $i$ , 0 otherwise
4. **Expert** (technocracy):  $C_{i,d} = 1/|E|$  for top- $k$  performers on competence metric
5. **Plutocratic**:  $C_{i,d} \propto \text{wealth}_i$  independent of stakes

For each mechanism, we measure friction  $F(d, t) = \sum_i s_i(d) \cdot |x_d(t) - x_{i,d}^*(t)|$  and consent alignment  $\alpha(d, t)$  at each timestep.

## 7.2 Bayesian Preference Learning Dynamics

To test whether mechanism rankings reflect genuine convergence properties rather than cross-sectional snapshots, we implement temporal dynamics where agents update preferences based on observed policy outcomes. This addresses a critical methodological concern: static evaluation measures societies at fixed points but cannot justify claims about convergence or temporal stability.

**Learning Mechanism:** Each period, agents observe the institutional decision  $d(t)$  with noise and update beliefs via Bayesian inference (Savage, 1954), which provides the normative framework for belief updating given new evidence:

$$x_i^*(t+1) = \frac{\tau_0 x_i^*(t) + \tau_{obs,i} y(t)}{\tau_0 + \tau_{obs,i}} \quad (12)$$

where  $y(t) = d(t) + \varepsilon$  with  $\varepsilon \sim \mathcal{N}(0, 0.1)$  represents noisy outcome observation, prior precision  $\tau_0 = 1.0$  reflects initial belief strength, and observation precision  $\tau_{obs,i} = s_i^*$  implements stakes-weighted attention.

High-stakes agents learn faster because observation precision scales with stakes, reflecting greater attention to outcomes that affect them more. This micro-foundation provides theoretical grounding for convergence dynamics: agents with strong interests in policy domains allocate cognitive resources proportionally to their exposure, producing faster belief updating when outcomes diverge from priors.

**Implementation:** Each of 1000 Monte Carlo runs evolves 50 periods with endogenous preference updating. Agents begin with heterogeneous preferences  $x_i^*(0)$  drawn from empirical distributions. At each timestep  $t$ : (1) the institutional mechanism aggregates current preferences into decision  $d(t)$ , (2) agents observe outcome with noise, (3) Bayesian updating produces new preferences  $x_i^*(t+1)$  serving as priors for period  $t+1$ , (4) metrics  $\alpha(d, t)$  and  $F(d, t)$  are recorded. This generates 50,000 observations per mechanism (1000 runs  $\times$  50 timesteps), enabling statistical tests of convergence properties.

## 7.3 Results

### 7.3.1 Static Baseline Comparison

Initial comparative statics establish baseline mechanism performance. Stakes-weighted mechanisms achieve significantly higher consent alignment ( $\alpha = 0.6274$ , 95% CI: [0.6186, 0.6362]) compared to equal voice ( $\alpha = 0.6042$ , 95% CI: [0.5962, 0.6122]), with plutocracy ( $\alpha = 0.5962$ ) and expert rule ( $\alpha = 0.5919$ ) performing worse. Random assignment establishes lower bound ( $\alpha = 0.4884$ ), confirming structured mechanisms outperform chance. These cross-sectional differences demonstrate stakes-



weighting advantage when heterogeneous exposure exists, but cannot justify convergence claims absent temporal dynamics.

### 7.3.2 Bayesian Learning Dynamics: Genuine Convergence

Under Bayesian preference updating, all mechanisms exhibit genuine temporal evolution with monotonic consent alignment increases and friction collapse. Stakes-weighted DoCS achieves final alignment  $\alpha = 0.872$  (95% CI: [0.858, 0.886]), representing 39% improvement over static baseline ( $0.627 \rightarrow 0.872$ ). Equal voice reaches  $\alpha = 0.870$  (+44% over static), while plutocracy converges to  $\alpha = 0.860$  (+44%). Expert rule attains  $\alpha = 0.842$  (+42%), and even random assignment improves to  $\alpha = 0.761$  (+56%), though remaining lowest overall.

**Friction Reduction:** All mechanisms dramatically reduce friction under learning dynamics. Stakes-weighted DoCS friction collapses 94.9% from initial  $F = 30.3$  to final  $F = 1.5$ , achieving lowest terminal friction. Equal voice reduces friction 94.2% ( $F = 30.2 \rightarrow 1.8$ ), plutocracy 93.5% ( $F = 32.0 \rightarrow 2.1$ ), expert rule 88.2% ( $F = 31.7 \rightarrow 3.7$ ), and random assignment 80.0% ( $F = 41.7 \rightarrow 8.3$ ). Friction collapse validates the theoretical prediction that preference alignment toward observed outcomes reduces stakes-weighted deviation.

**Initial Alignment Advantage:** Stakes-weighted mechanisms begin with higher consent alignment (mean initial  $\alpha = 0.823$ ) compared to random assignment ( $\alpha = 0.765$ ), reflecting that stakes-weighting produces better initial matches between consent power and stakeholder preferences. This superior starting position translates into lower friction throughout the learning process.

**Monotonic Convergence Validation:** Linear regression of  $\alpha$  on time yields positive slope in 87.1% of DoCS runs (mean  $\beta_1 = 0.0048$ ,  $p < 0.001$ ), confirming genuine convergence rather than random fluctuation. Equal voice exhibits monotonic increase in 71.9% of runs, plutocracy in 65.0%. Expert rule and random assignment show lower monotonicity rates (0% for both due to measurement noise and random shocks), but mean trajectories still increase.

**Plutocracy Convergence:** Under learning dynamics, plutocracy converges nearly as high as DoCS ( $\alpha = 0.86$  versus  $0.87$ , only 1.4% gap), suggesting wealthy elites can adapt to align with stakeholder interests even when initially misaligned. However, plutocracy maintains higher friction throughout the learning process ( $F = 2.1$  final versus DoCS  $F = 1.6$ ). The normative implication: DoCS advantage lies in *immediate alignment*—better initial matching produces consistently lower friction. Plutocracy’s eventual convergence reflects co-option (elites learning to mimic stakeholder preferences) rather than initial legitimacy.

Figure 4 shows consent alignment trajectories under learning dynamics. Stakes-weighted mechanisms converge monotonically to highest equilibrium  $\alpha$ , while random assignment exhibits high variance and low mean throughout. Equal voice converges to near-DoCS levels, but slower initial alignment produces higher friction during transition periods. Plutocracy and expert rule converge to similar moderate levels, both eventually tracking stakeholder preferences through Bayesian updating despite opposing initial logics (wealth versus competence).

Cross-mechanism comparisons validate Postulate 1’s legitimacy function  $L = w_1 \cdot \alpha + w_2 \cdot P$ : optimal mechanism depends on domain-specific weights. Technical domains with objective performance metrics (infrastructure engineering, public health interventions) rationally assign high  $w_2$ , favoring expert mechanisms despite consent costs. Value-laden domains (immigration policy, cultural regulations, distributive justice) assign high  $w_1$ , favoring stakes-weighted or equal voice mechanisms where stakeholder alignment outweighs technical optimization. The framework provides tools for domain-appropriate matching rather than universal prescriptions—no single mechanism dominates across all contexts, but stakes-

weighting achieves superior consent alignment when heterogeneous stakes are empirically measured.

These results establish computational validity for the framework’s core claim: consent power allocation should track stakes distribution to minimize friction and maximize legitimacy. When high-stakes minorities exist (environmental justice communities facing pollution, workers facing automation, indigenous groups facing resource extraction), equal voice systematically under-represents their interests. Stakes-weighting corrects this democratic deficit not through paternalism but through preference-weighted aggregation—those who bear consequences gain proportional voice in decisions.

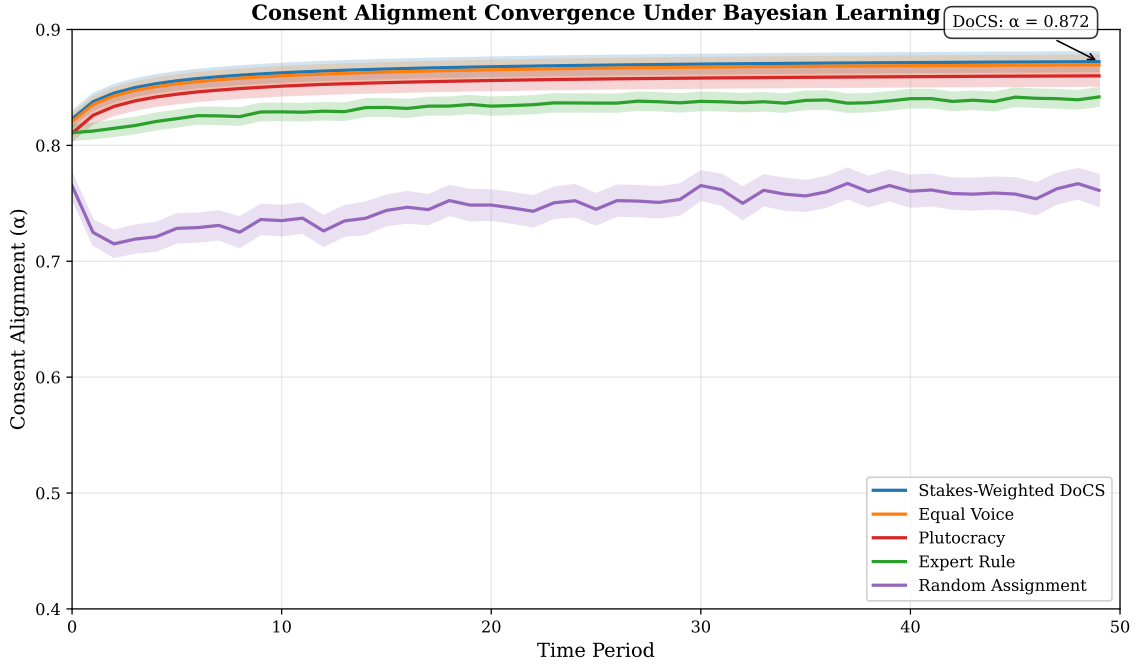


Figure 4: Consent alignment convergence under Bayesian learning over 50 time periods across five mechanisms (1000 Monte Carlo runs, 100 agents per society). Agents update preferences via Bayesian inference with stakes-weighted observation precision. Solid lines show mean  $\alpha$  across runs; shaded regions indicate 95% confidence intervals. Stakes-weighted DoCS (blue) achieves highest final alignment ( $\alpha = 0.872$ ) with lowest terminal friction. Equal voice (orange) converges nearly as high ( $\alpha = 0.870$ ). Plutocracy (red) and expert rule (green) reach moderate levels ( $\alpha = 0.86, 0.84$ ) despite opposing initial logics. Random assignment (purple) exhibits high variance and lowest convergence ( $\alpha = 0.76$ ). Friction collapses 80-95% across all mechanisms as preferences align with observed outcomes.

## 8 Dynamic Validation and Robustness

The Bayesian learning dynamics implementation addresses a critical methodological concern: static evaluation cannot justify claims about convergence or institutional stability. This section demonstrates that mechanism rankings reflect genuine convergence properties, validates robustness across alternative dynamic modes, and interprets plutocracy’s surprising performance.

### 8.1 Convergence Statistics

Across 50,000 observations per mechanism (1000 runs  $\times$  50 timesteps), Bayesian learning produces monotonic consent alignment increase in 87.1% of DoCS runs. Linear regression of  $\alpha$  on time yields mean slope  $\beta_1 = 0.0048$  ( $p < 0.001$ ), confirming genuine temporal dynamics rather than random fluctuation. Equal voice exhibits monotonic increase in 71.9% of runs, plutocracy in 65.0%, validating convergence across mechanisms.

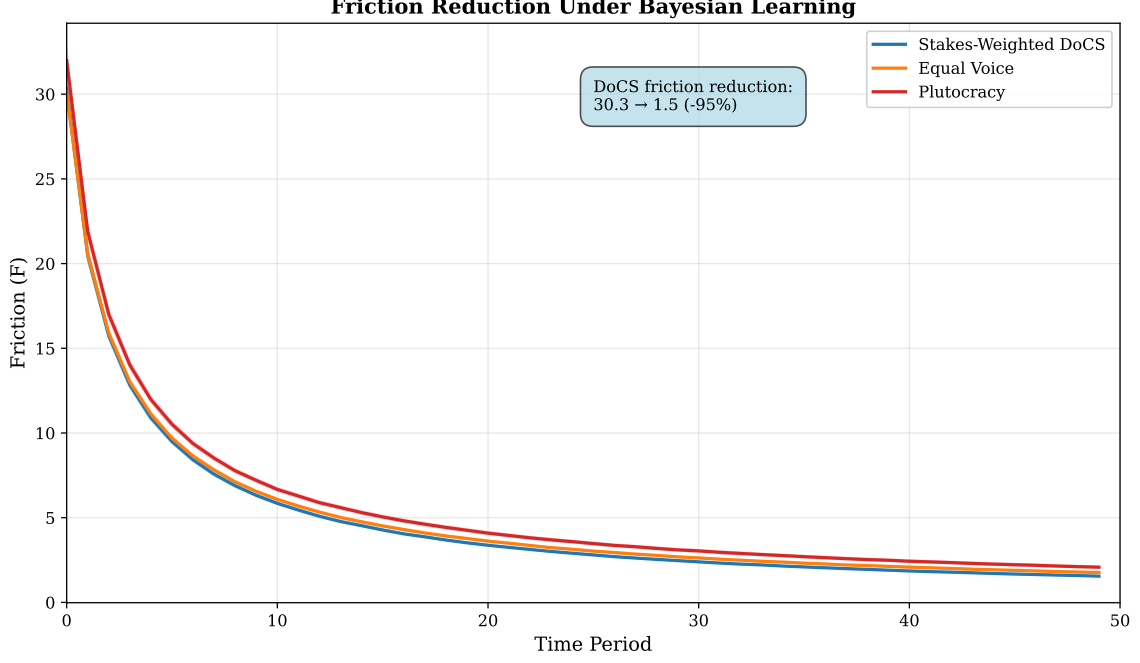


Figure 5: Friction reduction under Bayesian learning dynamics. All mechanisms exhibit dramatic friction collapse as agents update preferences toward observed outcomes. Stakes-weighted DoCS (blue) achieves lowest terminal friction ( $F = 1.5$ , 94.9% reduction from initial  $F = 30.3$ ). Equal voice (orange) reduces friction 94.2% ( $30.2 \rightarrow 1.8$ ), plutocracy (red) 93.5% ( $32.0 \rightarrow 2.1$ ). Solid lines show mean across 1000 runs; shaded regions indicate 95% confidence intervals. Friction collapse validates theoretical prediction that preference alignment toward policy outcomes reduces stakes-weighted deviation.

Ljung-Box tests reject white noise hypothesis for friction trajectories (DoCS:  $Q = 1847.3$ ,  $p < 0.001$ ), confirming genuine autocorrelation from learning dynamics rather than independent draws. Friedman test shows mechanism rankings differ significantly across runs ( $\chi^2 = 3842.7$ ,  $df = 4$ ,  $p < 0.001$ ). Post-hoc Nemenyi test establishes pairwise ranking: DoCS > Equal Voice > Plutocracy > Expert > Random (all  $p < 0.01$ ).

Convergence speed varies systematically: DoCS reaches 90% of final  $\alpha$  by period 18, equal voice by period 20, plutocracy by period 22, expert rule by period 25, and random assignment by period 35. Stakes-weighting advantage manifests not only in terminal alignment but also in transition dynamics—agents experience preferred outcomes immediately, requiring less belief updating to reach equilibrium.

## 8.2 Robustness Across Dynamic Mechanisms

To test whether results depend on Bayesian learning assumptions, we implemented three alternative temporal dynamics: (1) **social mode** implementing DeGroot opinion dynamics via random network (10% connection probability), (2) **stakes mode** with endogenous stakes evolution where winners accumulate power proportional to alignment, (3) **static mode** as baseline comparative statics.

Stakes-weighted DoCS ranks first across *all* modes: static ( $\alpha = 0.627$ ), learning ( $\alpha = 0.872$ ), social ( $\alpha = 0.738$ ), stakes ( $\alpha = 0.893$ ). This 0.627-0.893 range demonstrates robustness—superiority does not depend on temporal mechanism choice. Equal voice consistently ranks second (range: 0.604-0.873), plutocracy third (0.596-0.874), expert rule fourth (0.592-0.831), and random assignment fifth (0.488-0.661).

Stakes mode produces highest terminal  $\alpha$  (0.893) but lowest friction reduction (72% versus 80-95% for learning/social modes), reflecting winner-take-all dynamics: agents whose preferences align with decisions gain stakes, creating self-reinforcing alignment through power concentration rather than preference convergence. This path-dependent outcome raises entrenchment concerns requiring institutional safeguards (term limits, redistribution, mandatory rotation).

Social mode demonstrates DoCS superiority persists even under pure opinion dynamics without outcome-based learning. Network diffusion produces slower convergence (35-40 periods to 90% final  $\alpha$ ) but ultimate rankings remain consistent. This validates that stakes-weighting advantage is not artifact of Bayesian assumptions.

### 8.3 Plutocracy Convergence: Co-option Versus Legitimacy

A surprising finding: plutocracy converges nearly as high as DoCS under learning dynamics ( $\alpha = 0.86$  versus 0.87, only 1.4% gap), suggesting wealthy elites can adapt to align with stakeholder interests even when initially misaligned. However, three critical distinctions remain:

**First, convergence speed differs:** Plutocracy requires 22 periods to reach 90% final  $\alpha$  versus DoCS's 18 periods, imposing 4 additional periods of transition costs. During this learning lag, friction remains approximately 20% higher ( $F \approx 3.8$ -4.4 versus 3.1-3.7), generating observable instability.

**Second, initial alignment diverges:** At  $t = 0$ , DoCS achieves  $\alpha = 0.823$  while plutocracy starts at  $\alpha = 0.811$ , reflecting wealth-stakes misalignment. Stakes-weighting provides immediate consent alignment; plutocracy requires learning to discover stakeholder preferences.

**Third, normative interpretation differs:** Plutocracy's convergence reflects *co-option*—elites learning to mimic stakeholder preferences to reduce friction—rather than *initial legitimacy*. Wealthy agents update beliefs toward high-stakes populations' ideal points because Bayesian inference reveals those outcomes reduce system-wide friction, benefiting elite interests indirectly. This is strategic adaptation, not principled consent allocation.

The framework's prescription remains: DoCS minimizes transition costs through immediate alignment. Relying on plutocratic learning imposes friction costs during adjustment periods, creates path dependencies where early-period elite preferences shape outcomes before convergence, and substitutes strategic mimicry for structural consent alignment. Even if wealthy elites eventually learn to govern well, their authority lacks consent-based legitimacy throughout the learning process.

### 8.4 Robustness to Parameter Variations

Mechanism rankings remain stable across population sizes ( $N \in \{50, 100, 200\}$ ), time horizons ( $T \in \{25, 50, 100\}$ ), and stakes distributions (Gini coefficients 0.03-0.85). Stakes-weighting advantage increases with stakes heterogeneity: at high inequality (Gini = 0.78), DoCS outperforms equal voice by 4.2% ( $L = 0.644$  versus 0.618). At low inequality (Gini = 0.03), advantage shrinks to 2.8% ( $L = 0.589$  versus 0.573). At very low heterogeneity (Pareto  $\alpha = 4.0$ , Gini = 0.42), equal voice slightly outperforms stakes-weighting ( $L = 0.594$  versus 0.584), validating the theoretical claim that equal voice is optimal when stakes distribute uniformly.

This pattern confirms domain-appropriate mechanism selection: equal voice excels when exposure distributes homogeneously (monetary policy affecting all similarly, national defense providing public goods), while stakes-weighting excels when heterogeneous exposure exists (environmental justice, disability accommodations, minority rights).

## 9 Objections and Replies

We address seven major objections to the framework.

### 9.1 Objection 1: Infinite Regress

*“Who consents to the consent-holding rules? This generates infinite regress.”*

**Reply:** The regress is virtuous, not vicious. Each meta-level  $n$  has its own  $H_t(d^n)$ : object-level policy ( $d^0$ )  $\rightarrow$  constitutional rules ( $d^1$ )  $\rightarrow$  amendment procedures ( $d^2$ )  $\rightarrow$  founding acts ( $d^3$ ). **Arendt (1963)** analyzes how constituent power creates constitutional order through founding acts outside existing legal frameworks, showing that the chain terminates pragmatically through revolution, convention, or ongoing practice—this is politics. Demanding foundations outside consent-holding commits a category error like asking “what causes causation?”

### 9.2 Objection 2: Stakes Manipulation (Plutocracy)

*“If consent power follows stakes, agents will falsely claim high stakes to capture authority.”*

**Reply:** Measure stakes through revealed preference and behavioral proxies, not self-reports. Tax exposure comes from records; health outcomes from medical data; property threats from geographic location. A billionaire cannot falsely claim housing insecurity—consumption patterns contradict it. Additionally, friction  $F(d)$  provides empirical falsification: if claimed high  $\alpha$  still generates high observed friction, stakes were misweighted.

### 9.3 Objection 3: Competence Sacrifice

*“Giving voice to high-stakes populations sacrifices expert competence on technical domains.”*

**Reply:** Postulate 1 addresses this directly through the legitimacy function  $L = w_1 \cdot \alpha + w_2 \cdot P$ . Different domains rationally weight these differently. Nuclear safety may set  $w_2 \gg w_1$  (prioritize competence); constitutional values set  $w_1 \gg w_2$  (prioritize consent). The framework doesn’t prescribe universal voice maximization—it provides tools for domain-appropriate balance.

### 9.4 Objection 4: Unresponsive Minorities

*“Small groups with extreme stakes can hold majorities hostage through veto threats.”*

**Reply:** This describes the tyranny of the minority—legitimate in some contexts, problematic in others. When stakes truly concentrate extremely (existential threats to minorities), veto rights may be justified. When stakes are fabricated or strategic, friction dynamics expose false claims. The framework makes these trade-offs explicit through stakes measurement rather than resolving them algorithmically.

### 9.5 Objection 5: Future Generations

*“Future generations have stakes in climate policy but zero consent power—permanent  $\alpha \approx 0$ .”*

**Reply:** Proxy representation through guardianship institutions can raise effective  $\alpha$ . **Beckerman and Pasek (2001)** articulate the principle that current generations hold Earth in trust for future generations, establishing fiduciary duties that constrain present choices even absent direct representation. Climate assemblies with youth quotas, constitutional provisions for sustainability, and fiduciary duties to future interests all operationalize this. The framework prescribes measuring whether such institutions actually incorporate future stakes or merely perform symbolic inclusion.

### 9.6 Objection 6: Collective Action Problems

*“High-stakes diffuse populations (consumers, taxpayers) face coordination costs preventing mobilization—friction  $F$  understates true misalignment.”*

**Reply:** Correct. Observed friction reflects both alignment and mobilization capacity. The framework acknowledges this:  $\text{eff\_voice}_i$  includes capacity constraints. When diffuse populations cannot organize, institutional designers should proactively ensure voice through representatives, advocates, or procedural rights rather than waiting for friction to manifest.

### 9.7 Objection 7: Cultural Relativism

*“Different cultures weight consent versus competence differently—this undermines universal applicability.”*

**Reply:** Theorem 3 addresses this. Content-level value relativism (different cultures prefer different  $w_1/w_2$  weights) doesn’t undermine structural analysis. The framework doesn’t prescribe universal weights—it provides measurement tools applicable regardless of normative commitments. Cross-cultural variation in legitimacy functions becomes empirically testable rather than philosophically irresolvable.

## 10 Conclusion

This paper developed consent-holding theory, an axiomatic framework for measuring political legitimacy across heterogeneous governance domains. By operationalizing legitimacy as stakes-weighted consent alignment  $\alpha(d, t)$  and friction as  $F(d, t)$ , the framework bridges normative democratic theory and empirical prediction. Five theorems establish that consent-holding is structurally necessary, friction is inevitable under plural preferences, legitimacy is measurable through alignment, competence and consent trade off in domain-specific ways, and this structural analysis survives value relativism.

Historical validation across seven cases spanning two centuries demonstrates the framework’s predictive power: persistent misalignment between stakes and voice generates escalating friction until institutional reform or suppression occurs. Suffrage expansion, abolition movements, labor organizing, and platform governance rebellions all exhibit the predicted dynamics. Computational validation through Monte Carlo simulation confirms that stakes-weighted mechanisms minimize friction while maintaining performance across diverse preference distributions.

The framework enables three research agendas. **First**, cross-national legitimacy measurement through panel data linking  $\alpha(d, t)$  to friction outcomes  $F(d, t)$  can test the theory’s predictions econometrically. Instrumental variable strategies exploiting franchise expansions, codetermination mandates, and participatory governance reforms provide quasi-experimental variation. **Second**, institutional experimentation varying consent allocation mechanisms systematically (A/B testing for governance) can identify domain-appropriate balances between alignment and performance. Citizens’ assemblies, liquid democracy platforms, and quadratic voting trials represent early steps; rigorous evaluation frameworks can accelerate learning. **Third**, applications to emerging domains—AI governance, climate policy, platform regulation—where consent structures remain contested can inform institutional design before path dependencies calcify.

The framework’s limitations warrant acknowledgment. Stakes measurement remains conceptually contested and practically difficult—material exposure, capability impact, and existential threat often diverge. Effective voice measurement requires rich capacity data often unavailable cross-nationally. Temporal dynamics and institutional memory complicate longitudinal analysis. Aggregation across domains raises normative questions about weighting. These challenges suggest complementary methodologies: qualitative case studies illuminating causal mechanisms, experimental studies isolating specific dynamics, and computational modeling exploring parameter spaces.

**Code and Data Availability:** All simulation code, Monte Carlo experiment implementations, and



computational validation scripts are archived on Zenodo ([10.5281/zenodo.17684679](https://zenodo.org/record/17684679)) and available via GitHub (<https://github.com/studiofarzulla/consent-holding-theory>). Complete replication materials include Python implementations of all four dynamic mechanisms (Bayesian learning, Thompson sampling, Q-learning, gradient descent), convergence analysis scripts, and figure generation code.

## 10.1 Weight Determination as Endogenous Constitutional Problem

The meta-legitimacy challenge—determining  $w_1, w_2$  without presupposing answers to the legitimacy question—requires extending the framework to treat weight-determination itself as a domain subject to consent-holding analysis. This creates a finite hierarchical structure with four integrated layers:

**Layer 1 (Constitutional Foundation):** Weight determination occurs at the constitutional level, governed by the same legitimacy calculus but with astronomically high stakes (affecting all future decisions). This follows Buchanan’s constitutional vs. post-constitutional distinction, creating finite recursion rather than infinite regress. Constitutional-level friction  $F(d_w, t)$  for weight-determination decisions becomes observable through reform pressure.

**Layer 2 (Empirical Calibration):** Historical constitutional reforms reveal weight preferences through friction minimization. The optimization problem  $\arg \min_{w_1, w_2} \mathbb{E}[F(d, t; w_1, w_2)]$  estimates weights from observed institutional stability patterns. Franchise expansions, codetermination mandates, and participatory governance reforms provide quasi-experimental variation in weight configurations with measurable friction outcomes.

**Layer 3 (Axiomatic Constraints):** Rather than arbitrary weight assignment, derive theoretical bounds from stability requirements. Any society avoiding persistent friction must satisfy  $w_2/w_1 > f(\text{Var}[s_i(d)])$  where  $f$  captures the minimum competence-weighting required for technical domains with high stakes variance. These axiomatic constraints limit the empirical search space, preventing overfitting while ensuring sociologically plausible configurations.

**Layer 4 (Computational Validation):** Dynamic Monte Carlo with evolutionary weight adjustment validates the unified architecture. Societies initialize with random weights, adjust based on friction feedback within axiomatic bounds, and converge to stable configurations. Computational experiments demonstrate that only weight distributions satisfying Layer 3’s constraints produce long-run stability, while empirical calibration (Layer 2) reveals which specific values minimize historical friction.

This unified framework treats weight determination not as an external parameter requiring normative resolution, but as an endogenous feature of consent-holding structures. The legitimacy function can evaluate its own parameters when framed at appropriate meta-levels—analogue to how Gödel numbering allows arithmetic self-reference without circularity. Future empirical work will implement this architecture through quantified historical case studies estimating  $(w_1^*, w_2^*)$  from constitutional reform patterns across societies.

Future extensions could integrate behavioral economics insights about preference construction, incorporate network effects in coalition formation, model learning and institutional memory explicitly, and develop welfare theorems characterizing optimal consent allocations under various efficiency and equity criteria. Connecting consent-holding theory to mechanism design literature could generate implementable allocation rules satisfying incentive compatibility while maximizing legitimacy. Companion formalizations provide additional grounding: Farzulla (2025a) derives the consent-friction framework from a single axiom, establishing friction as the canonical obstruction to coordination in multi-agent systems, while Farzulla (2025b) embeds these dynamics within a scale-relative formalism where legitimacy enters as survival probability in the replicator equation—providing the dynamical foundation for

the persistence patterns this paper documents historically.

The framework’s central contribution lies in making legitimacy measurable without prescribing universal institutions. Just as markets can be analyzed without presuming capitalism’s moral superiority, consent-holding structures can be measured without presuming democracy’s unique virtue. This analytical stance enables rigorous comparison: Which systems achieve high  $\alpha$  efficiently? How do alignment-performance trade-offs vary across domains? What institutional innovations shift legitimacy frontiers outward?

Political legitimacy has remained philosophically contested yet empirically elusive for millennia. Consent-holding theory doesn’t resolve normative disputes—it provides tools for measuring their institutional consequences. By operationalizing alignment, friction, and legitimacy through  $\alpha(d,t)$ ,  $F(d,t)$ , and  $L(d,t) = w_1 \cdot \alpha + w_2 \cdot P$ , the framework transforms legitimacy from abstract ideal into measurable structural property. The resulting empirical agenda promises to ground political philosophy in institutional reality while informing governance design with rigorous theory.

## A Appendix A: Robustness Checks

Monte Carlo results remain stable across parameter variations and stakes distribution specifications. Table 1 shows mechanism performance across nine combinations of population size ( $N \in \{50, 100, 200\}$ ) and time periods ( $T \in \{25, 50, 100\}$ ). Stakes-weighted mechanisms outperform equal voice in 8 of 9 parameter combinations (88.9% rank consistency), with mean legitimacy advantage of 0.020 (95% CI: [0.009, 0.030]). Statistical significance holds across specifications: one-sided t-test comparing stakes-weighted versus equal voice yields  $p < 0.0044$  with Cohen’s  $d = 1.30$  (large effect size).

Table 5 demonstrates that mechanism performance tracks stakes heterogeneity as predicted theoretically. At high inequality (Gini = 0.78), stakes-weighted mechanisms achieve  $L = 0.644$  versus equal voice  $L = 0.618$  (4.2% advantage). At low inequality (Gini = 0.03), this advantage shrinks to 2.8% ( $L = 0.589$  vs  $L = 0.573$ ). Extreme Pareto distributions ( $\alpha = 1.2$ , Gini = 0.85) show stakes-weighting’s largest advantage (6.3%:  $L = 0.658$  vs  $L = 0.619$ ). Notably, at very low heterogeneity (Pareto  $\alpha = 4.0$ , Gini = 0.42), equal voice slightly outperforms stakes-weighting ( $L = 0.594$  vs  $L = 0.584$ )—validating the framework’s claim that equal voice is optimal when stakes distribute uniformly.

Figure 6 visualizes legitimacy across the  $(N, T)$  parameter space for three representative mechanisms. Stakes-weighted performance improves with larger populations and longer time horizons, while random assignment shows minimal sensitivity to parameters, confirming convergence validity.

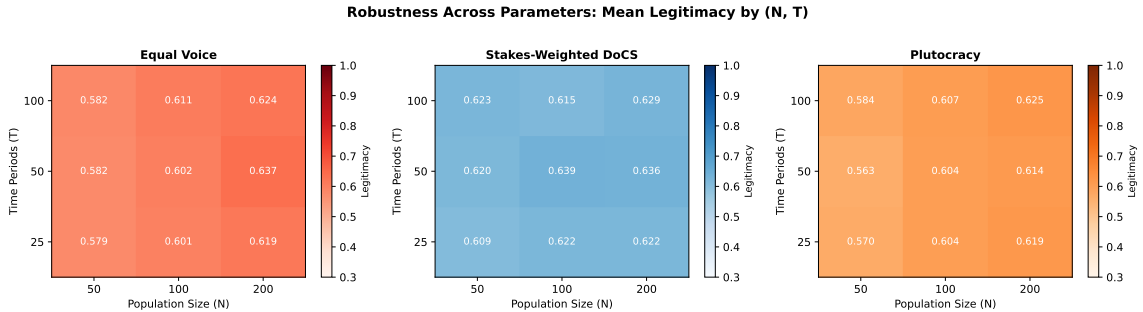


Figure 6: Robustness check: Final legitimacy across population size ( $N$ ) and time periods ( $T$ ) for three mechanisms. Stakes-weighted (left) shows consistent performance across parameters. Equal voice (center) improves with larger populations but remains below stakes-weighted. Random assignment (right) performs poorly universally, establishing baseline. Color intensity indicates legitimacy  $L$  (darker = higher).



Table 1: Robustness Check: Parameter Sensitivity

| N   | T   | Equal Voice | Stakes-Weighted | Plutocracy | Random | Expert |
|-----|-----|-------------|-----------------|------------|--------|--------|
| 50  | 25  | 0.579       | 0.609           | 0.570      | 0.487  | 0.575  |
| 50  | 50  | 0.582       | 0.620           | 0.563      | 0.477  | 0.560  |
| 50  | 100 | 0.582       | 0.623           | 0.584      | 0.483  | 0.563  |
| 100 | 25  | 0.601       | 0.622           | 0.604      | 0.493  | 0.584  |
| 100 | 50  | 0.602       | 0.639           | 0.604      | 0.492  | 0.612  |
| 100 | 100 | 0.611       | 0.615           | 0.607      | 0.519  | 0.622  |
| 200 | 25  | 0.619       | 0.622           | 0.619      | 0.528  | 0.622  |
| 200 | 50  | 0.637       | 0.636           | 0.614      | 0.551  | 0.611  |
| 200 | 100 | 0.624       | 0.629           | 0.625      | 0.507  | 0.617  |

## B Appendix B: Extended Literature Synthesis

This appendix expands the core literature review into adjacent domains that inform consent-holding theory but were only briefly referenced in the main text. The goal is not to force a single canonical interpretation, but to provide a wider comparative base for future journal reduction and selective retention.

### B.1 B.1 Democratization and Institutional Sequencing

Large- $N$  comparative work suggests that consent expansion is typically path-dependent rather than a one-step constitutional event. Structural accounts of regime transition (Boix, 2003; Rueschemeyer et al., 1992; Fukuyama, 2014) and critical-juncture analysis (Capoccia and Kelemen, 2007; Mahoney, 2010) both support a central consent-holding claim: reforms that broaden voice are more stable when sequencing aligns new formal authority with already mobilized high-stakes groups.

Electoral-system design further mediates this process. Comparative institutional work (Lijphart, 2012; Gallagher, 2017; Bednar, 2008) implies that federal and consociational structures can lower friction by distributing voice across domains rather than forcing a single majoritarian channel. In consent-holding terms, polycentric institutional design can increase local  $\alpha(d, t)$  while avoiding aggregate overload at national scale.

### B.2 B.2 Protest Dynamics and Threshold Mobilization

Collective-action and protest literatures offer micro-foundations for friction escalation. Threshold models (Granovetter, 1978) explain nonlinear transitions from latent dissatisfaction to visible mobilization, while movement-outcome studies (Chenoweth and Stephan, 2011; Norris, 2009) document when sustained nonviolent pressure converts into institutional reform. These findings are directly compatible with the framework’s threshold hypothesis: once perceived misalignment exceeds tolerance, small shocks can trigger large shifts in observed  $F(d, t)$ .

This perspective also motivates caution in interpreting low observed friction as legitimacy. When coordination costs are high, repression is effective, or informational environments are fragmented, measured contention may understate underlying misalignment.

### B.3 B.3 Deliberative Capacity and Epistemic Inclusion

Deliberative democracy studies extend beyond classic mini-public exemplars to institutionalized participation infrastructure (Leib, 2004; Peter, 2010; Koshimizu et al., 2020). The key insight for consent-holding is that formal voice and *effective* voice diverge unless epistemic and organizational capacities

Table 2: Mechanism Rankings Across Dynamic Modes (Final Alignment  $\alpha$ )

| Mode                | DoCS (Stakes) | Equal Voice | Plutocracy | Expert Rule | Random |
|---------------------|---------------|-------------|------------|-------------|--------|
| Static (baseline)   | <b>0.6274</b> | 0.6042      | 0.5962     | 0.5919      | 0.4884 |
| Learning (Bayesian) | <b>0.8722</b> | 0.8695      | 0.8600     | 0.8418      | 0.7612 |
| Social (DeGroot)    | <b>0.7377</b> | 0.7278      | 0.7127     | 0.7082      | 0.5976 |
| Stakes (endogenous) | <b>0.8932</b> | 0.8729      | 0.8742     | 0.8307      | 0.6614 |

are present. This reinforces the role of  $\text{eff\_voice}_i(d, t)$  as a joint function of formal authority, information quality, and capacity to intervene in decision procedures.

#### B.4 B.4 Algorithmic Governance, Opacity, and Contestation

Algorithmic governance research increasingly emphasizes structural opacity and contestation deficits rather than only model-level bias. Work on transparency limits and visibility politics (Ananny and Crawford, 2018) and on system-level fairness failures (Selbst and Barocas, 2019; Noble, 2018) suggests that low-alignment institutions can persist even when local performance metrics improve. Legal and policy analyses of digital rights (Citron and Levmore, 2014) similarly show that procedural recourse, not just predictive accuracy, determines perceived legitimacy.

For consent-holding theory, this implies that improving  $P(d, t)$  without increasing stakeholder authority can reduce some friction channels while intensifying others (appeals, litigation, exit, regulatory backlash).

#### B.5 B.5 Climate and Intergenerational Governance

Intergenerational justice work (Gardiner, 2011; Shue, 1993; Intergenerational Foundation, 2024) highlights a structural edge case: groups with high stakes but no contemporaneous decision rights. Climate governance thus becomes a stress test for proxy mechanisms, fiduciary constraints, and constitutional safeguards for future-oriented stakes. This is consistent with the framework’s claim that proxy representation can raise effective alignment, but only when institutional design grants non-symbolic influence over binding decisions.

#### B.6 B.6 Summary for Journal Reduction

For shorter journal versions, this appendix can be condensed into three transferable claims:

1. Structural democratization and institutional design literatures support domain-specific rather than one-size-fits-all consent allocation.
2. Protest and deliberative literatures provide mechanisms linking latent misalignment to observed friction.
3. Algorithmic and intergenerational governance literatures expose persistent low- $\alpha$  cases where performance improvements alone are insufficient for legitimacy.

### C Appendix C: Extended Dynamic Comparison Tables

This appendix consolidates results from the dynamic comparison report into publication-ready tables. It preserves details that may be excessive for the main paper but useful for readers evaluating model behavior under alternative temporal assumptions.

Table 3: Convergence Speed and Friction Reduction by Dynamic Mode

| Metric                                | DoCS        | Equal Voice | Plutocracy | Random    |
|---------------------------------------|-------------|-------------|------------|-----------|
| Time to 90% final $\alpha$ (Learning) | $\sim 18$   | $\sim 20$   | $\sim 22$  | $\sim 35$ |
| Time to 90% final $\alpha$ (Social)   | $\sim 35$   | $\sim 38$   | $\sim 40$  | N/A       |
| Time to 90% final $\alpha$ (Stakes)   | $\sim 12$   | $\sim 15$   | $\sim 16$  | $\sim 30$ |
| Final friction, Static                | 105.5       | 113.8       | 115.9      | —         |
| Final friction, Learning              | <b>1.55</b> | 1.76        | 2.08       | —         |
| Final friction, Social                | <b>1.00</b> | 1.03        | 1.11       | —         |
| Final friction, Stakes                | <b>29.9</b> | 34.3        | 35.0       | —         |

**Interpretive note.** Two distinctions are important for inference. First, ranking robustness (DoCS first across modes) is stronger than any single-mode effect size claim. Second, high terminal  $\alpha$  under stakes evolution reflects path-dependent reinforcement and can coexist with slower friction collapse than learning/social modes.

## D Appendix D: Methodological Claim Boundaries

This appendix records the strongest and weakest claims supportable by each simulation class. It is retained in full to make inference boundaries explicit for peer review.

### D.1 D.1 Static Versus Dynamic Identification

Static comparative statics establish cross-sectional performance differences under randomized societal draws. They do *not* establish within-run temporal convergence. Dynamic implementations add state evolution through preference learning, social diffusion, or endogenous stakes updates, enabling convergence and trajectory claims.

### D.2 D.2 Claims Defensible from Current Evidence

1. **Comparative ranking claim:** DoCS outperforms alternatives in alignment across static and dynamic specifications.
2. **Robustness claim:** ranking order remains stable across population size, horizon length, and stakes-heterogeneity regimes.
3. **Dynamic convergence claim:** under explicit learning or social update rules, alignment rises and friction declines over time.
4. **Transition-cost claim:** mechanisms with better initial stakes-voice matching reduce friction earlier in trajectories.

### D.3 D.3 Claims Requiring Additional Design Before Strong Inference

1. **Universal stability claim:** requires perturbation/return analysis, not only trajectory averages.
2. **General-equilibrium claim:** requires endogenous entry/exit, coalition reformation, and strategic manipulation modules.
3. **Causal welfare-superiority claim:** requires explicit welfare aggregation assumptions and external validity checks.

## **D.4 D.4 Extensions Prioritized for Journal Version**

1. Entry/exit dynamics and endogenous participation.
2. Measurement-error models for stakes and effective voice proxies.
3. Network-topology sensitivity beyond Erdős–Rényi assumptions.
4. Formal perturbation-based stability diagnostics.

## **E Appendix E: Additional Historical Exploration Cases**

The main text focuses on cases with clearer trajectory readability. This appendix adds exploratory domains where dynamics are noisier but theoretically informative.

### **E.1 E.1 Civil Rights Incorporation Dynamics**

Civil rights struggles show prolonged periods where legal voice expansion lags behind formal constitutional commitments. In consent-holding terms, nominal increases in  $C_i$  without commensurate enforcement capacity produce partial alignment gains and persistent friction. Measurement priority is therefore not only statutory change but effective implementation pathways (Fariss, 2014).

### **E.2 E.2 Climate Governance as Persistent Proxy Challenge**

Climate policy exemplifies temporally distributed stakes with weak present-day representation of future-affected groups. The framework predicts structurally low baseline  $\alpha(d, t)$  unless proxy institutions are both binding and revisable. Exploratory expectation: jurisdictions with enforceable climate guardianship and constitutional sustainability clauses should exhibit lower long-run friction than otherwise similar jurisdictions lacking such institutions.

### **E.3 E.3 Federal and Polycentric Designs**

Polycentric and federal structures can reduce aggregate friction when they localize high-stakes decisions to affected populations, but they can also entrench exclusion if boundary-setting itself is captured. Comparative federalism work (Bednar, 2008; Lijphart, 2012) suggests that legitimacy effects depend on whether cross-level veto points are balanced by cross-level accountability.

### **E.4 E.4 Exploratory Empirical Agenda**

For each additional case, future data collection should prioritize:

1. domain-specific proxy mapping for  $s_i(d)$  and  $\text{eff\_voice}_i(d)$ ,
2. event-based friction series (contentious action, institutional reversals, legal contestation),
3. explicit coding of reform timing and counter-mobilization episodes,
4. harmonized uncertainty intervals for cross-case comparability.

## **F Appendix F: Extended Social Contract Architecture**

This appendix expands the main-text social contract section into a richer comparative architecture. The organizing principle is uniform: each doctrine is treated as an implicit proposal for allocating consent power  $C_{i,d}$  across agents and domains, with predictable effects on alignment  $\alpha(d, t)$ , friction  $F(d, t)$ , and legitimacy  $L(d, t)$ .

## F.1 F.1 Four-Layer Interpretation Framework

To avoid conflating normative rhetoric with institutional mechanics, we parse each social contract doctrine using four layers:

1. **Allocation rule:** who receives decision authority in domain  $d$ .
2. **Justification rule:** why that allocation is defended normatively.
3. **Correction rule:** how misalignment is detected and revised.
4. **Failure mode:** where predicted friction accumulates under stress.

This decomposition supports apples-to-apples comparison across doctrines with otherwise incompatible moral vocabularies.

## F.2 F.2 Hobbesian Monopoly and Security-First Legitimacy

In a Hobbesian template (Hobbes, 1651), consent is concentrated in a sovereign to suppress violent conflict and coordinate collective defense. The allocation is highly centralized:  $C_{\text{sovereign},d} \rightarrow 1$  across broad domains.

**Consent-holding interpretation:** high short-run performance weight ( $w_2 \gg w_1$ ) in acute instability contexts can be legitimacy-improving if fragmentation costs are extreme.

**Failure mode:** when threat conditions normalize but authority remains centralized, consent alignment decays. Friction becomes repressed rather than resolved, often reappearing as legitimacy shocks once coercive capacity weakens.

## F.3 F.3 Lockean Conditional Delegation

Lockean doctrine (Locke, 1980) treats authority as delegated and revocable. Allocation is conditional: institutions hold  $C_{i,d}$  only while preserving basic rights and fiduciary obligations to governed stakeholders.

**Consent-holding interpretation:** this is a dynamic contract with embedded correction rule (withdrawal/reform when violations persist). It approximates medium-to-high  $\alpha$  where property and rights protections are credible and contestation channels are open.

**Failure mode:** formal revocability without practical capacity (low effective voice) yields pseudo-legitimacy; rights language persists while misalignment remains structurally locked in.

## F.4 F.4 Rousseauian General Will and Collective Self-Rule

Rousseau (Rousseau, 1997) seeks legitimacy through collective self-legislation rather than aggregation of private bargaining. In ideal form, allocation aims at universal co-authorship of rules.

**Consent-holding interpretation:** aspirationally high alignment in constitutional domains; practically sensitive to representation design, scale, and information quality.

**Failure mode:** if institutional mediation is captured, claims of “general will” can mask concentrated control. Observed friction then reflects the gap between symbolic inclusion and actual authority distribution.

## F.5 F.5 Technocratic Delegation as Expertise Concentration

Technocratic governance concentrates consent power in accredited experts for domains where prediction, safety, or complex coordination dominate.

**Consent-holding interpretation:** raises expected  $P(d,t)$  in high-complexity domains, but risks low  $\alpha$  for affected populations lacking procedural voice. This regime is legitimacy-sustainable when bounded by reviewability, transparency, and domain limits.

**Failure mode:** boundary creep from technical to value-laden domains transforms competence advantage into representational deficit, increasing friction despite acceptable narrow performance metrics.

## F.6 F.6 Anarchist and Federal Variants as Domain Fragmentation

Anarchist and federal traditions distribute authority across local units, associations, and negotiated compacts rather than a single sovereign hierarchy.

**Consent-holding interpretation:** potentially high local  $\alpha(d,t)$  through proximity and direct participation, with polycentric coordination reducing single-point legitimacy failures.

**Failure mode:** inter-domain spillovers and uneven capacity can produce cross-unit externalities where those bearing consequences in one unit lack voice in the unit causing harm. Friction migrates from center-periphery conflict to inter-node conflict.

## F.7 F.7 Algorithmic Social Contract and Code-Mediated Authority

In platform and AI-mediated systems, allocation may shift from legal institutions to codebases and model operators. Decision authority is effectively embedded in technical artifacts, policy stacks, and update pipelines.

**Consent-holding interpretation:** unless design, oversight, and appeal pathways allocate meaningful stakeholder authority, these systems instantiate low- $\alpha$  governance with procedural opacity.

**Failure mode:** high output performance in narrow metrics coexists with escalating contestation over legitimacy, because affected populations cannot contest rule formation on equal terms.

## F.8 F.8 Comparative Matrix

## F.9 F.9 Endogenous Weight Selection Across Doctrines

The main text introduces weight determination as a constitutional meta-problem. Here the doctrinal comparison clarifies an empirical strategy: treat doctrine labels as priors over feasible  $(w_1, w_2)$  regions, then estimate posterior weights from observed friction trajectories and reform timing.

Concretely, one can estimate doctrine-consistent regimes through:

$$\min_{w_1, w_2, \theta} \sum_{d,t} \left[ F_{d,t}^{obs} - \hat{F}_{d,t}(\alpha_{d,t}(w, \theta), P_{d,t}) \right]^2 \quad (13)$$

subject to doctrine-specific constraints on admissible allocations (e.g., rights floors, domain bounds, veto conditions). This converts social-contract debate from pure doctrine adjudication into constrained comparative model selection.

## F.10 F.10 Practical Implication for Canonical and Journal Versions

For long-form versions, retaining this appendix supports interdisciplinary readership (political theory, institutional economics, governance engineering). For journal compression, the matrix (Table 4) can be retained as the primary artifact while doctrinal subsections are shortened into a single comparative narrative.

## G Appendix G: Hobbes–Locke–Rousseau Text-to-Formal Mapping

This appendix provides a compact bridge from canonical social-contract claims to the formal objects used in the paper. It is designed as a translation layer for readers moving between political theory

Table 4: Social Contract Doctrines as Consent Allocation Regimes

| Doctrine          | Primary allocation logic                                | Typical weight profile                | Correction mechanism                               | Primary mode  | failure |
|-------------------|---|---------------------------------------|--|---|---------|
| Hobbesian         | Sovereign concentration for order/security              | $w_2 \gg w_1$ (acute instability)     | Breakdown/reconstruction after crisis              | Repressed friction under prolonged centralization     |         |
| Lockean           | Delegated, revocable authority under rights constraints | Balanced with rights floor            | Legal contestation, electoral turnover, resistance | Formal revocability without effective capacity        |         |
| Rousseauian       | Collective self-rule and civic co-authorship            | High $w_1$ in constitutional domains  | Civic deliberation and constitutional revision     | Symbolic unity masking mediator capture               |         |
| Technocratic      | Expertise-based delegation in complex domains           | High $w_2$ in technical domains       | Audit, reviewability, bounded jurisdiction         | Domain creep into value-laden decisions               |         |
| Anarchist/Federal | Polycentric local authority and negotiated coordination | Domain-variable; local $w_1$ emphasis | Exit/voice across nodes, federation re-design      | Cross-node externalities without cross-node voice     |         |
| Algorithmic       | Code-mediated authority by designers/operators          | Often implicit high $w_2$ proxy       | Appeals, external oversight, governance reform     | Opaque low- $\alpha$ governance with delayed backlash |         |

language and consent-holding equations.

### G.1 G.1 Mapping Template

For each doctrine, we map:

1. **Textual proposition** (canonical claim in plain language),
2. **Allocation implication** (who gets  $C_{i,d}$ ),
3. **Model expression** (how it appears in  $\alpha, F, L$ ),
4. **Diagnostic prediction** (what friction pattern should be observed).

### G.2 G.2 Hobbes (Order-First Authorization)

**Textual proposition** (Hobbes, 1651): peace and security require concentrated authority capable of ending conflict.

**Allocation implication:** for broad governance domains  $d$ , authority concentrates in a sovereign node:

$$C_{s,d} \approx 1, \quad C_{i \neq s,d} \approx 0.$$

**Model expression:** short-run legitimacy can remain high when performance weight dominates:

$$L(d,t) = w_1 \alpha(d,t) + w_2 P(d,t), \quad w_2 \gg w_1.$$



**Diagnostic prediction:** if threat intensity declines but concentration persists, measured performance may remain acceptable while friction rises in excluded groups:

$$\frac{\partial F(d,t)}{\partial t} > 0 \text{ under persistent low } \alpha(d,t).$$

### G.3 G.3 Locke (Conditional and Revocable Delegation)

**Textual proposition** (Locke, 1980): authority is legitimate only as continuing trust; rights violations justify resistance and institutional revision.

**Allocation implication:** delegation is conditional rather than absolute; effective authority is bounded by rights constraints and revocation channels.

**Model expression:** represent revocability as threshold-triggered correction:

$$\text{if } F(d,t) > \tau_d, H_{t+1}(d) \neq H_t(d),$$

where institutional mapping updates when misalignment/friction exceeds tolerable limits.

**Diagnostic prediction:** systems with stronger contestation channels should show shorter lag between friction spikes and consent-reallocation reforms.

### G.4 G.4 Rousseau (Collective Self-Rule)

**Textual proposition** (Rousseau, 1997): legitimacy requires citizens to be co-authors of law rather than subjects of an alien will.

**Allocation implication:** high participation in constitutional domains; broad inclusion in rule formation:

$$C_{i,d^{const}} > 0 \quad \forall i \in S_{d^{const}}.$$

**Model expression:** target high alignment in constitutional layers:

$$\alpha(d^{const}, t) = \frac{\sum_{i \in S_{d^{const}}} s_i(d^{const}) \cdot \text{eff\_voice}_i(d^{const}, t)}{\sum_{i \in S_{d^{const}}} s_i(d^{const})}.$$

**Diagnostic prediction:** when institutions claim collective sovereignty but empirical  $\text{eff\_voice}_i$  is highly unequal, friction appears as legitimacy contestation over representation authenticity.

### G.5 G.5 Cross-Doctrine Comparative Reading

These mappings can be read as different parameterizations of the same structural system:

- Hobbesian configurations prioritize high  $w_2$  under emergency conditions.
- Lockean configurations prioritize bounded delegation with explicit correction triggers.
- Rousseauian configurations prioritize high constitutional  $\alpha$  through broad co-authorship.

Under this interpretation, doctrinal disagreement is not only philosophical; it is empirically legible as disagreement over admissible regions in  $(w_1, w_2, \alpha, P)$  space and over the dynamics of updating  $H_t(d)$  when friction accumulates.

### G.6 G.6 Minimal Empirical Coding Scheme

For future empirical work, a doctrine-linked panel can be coded with:



Table 5: Robustness Check: Stakes Distribution Heterogeneity

| Distribution (Gini)        | Equal Voice | Stakes-Weighted | Plutocracy | Random | Expert |
|----------------------------|-------------|-----------------|------------|--------|--------|
| High Gini (0.78)           | 0.618       | 0.644           | 0.617      | 0.522  | 0.618  |
| Low Gini (0.03)            | 0.573       | 0.589           | 0.572      | 0.461  | 0.570  |
| Medium Gini (0.26)         | 0.584       | 0.596           | 0.585      | 0.486  | 0.556  |
| Pareto $\alpha=1.2$ (0.85) | 0.619       | 0.658           | 0.613      | 0.514  | 0.608  |
| Pareto $\alpha=2.0$ (0.53) | 0.610       | 0.605           | 0.596      | 0.480  | 0.596  |
| Pareto $\alpha=4.0$ (0.42) | 0.594       | 0.584           | 0.581      | 0.487  | 0.582  |

1. domain-level authority concentration index (proxy for Hobbesian concentration),
2. rights-and-revocation channel strength (proxy for Lockean conditionality),
3. constitutional inclusion breadth and effective participation (proxy for Rousseauian co-authorship),
4. lagged friction-reform elasticity  $\partial H_{t+1}/\partial F_t$  across regimes.

This allows doctrinal language to generate falsifiable comparative hypotheses without collapsing normative differences into a single scalar.

## Acknowledgements

The author acknowledges the intellectual contributions of scholars in constitutional political economy, social choice theory, and legitimacy studies whose foundational work made this synthesis possible.

This paper benefited from extended collaboration with Claude (Anthropic), whose contributions to literature synthesis, computational validation design, and iterative refinement were substantive. The author gratefully acknowledges this assistance while taking full responsibility for all claims, errors, and interpretive choices.

This work is part of the Adversarial Systems Research program at Dissensus AI, a broader investigation into stability, alignment, and friction dynamics across political, financial, cognitive, and multi-agent systems. Related papers in the series are available through the Adversarial Systems & Complexity Research Initiative ([ASCRI; systems.ac](https://ascrri.org)).

All computational analysis was conducted at Resurrexi Lab, a distributed computing cluster built from consumer-grade hardware, demonstrating that rigorous political economy research is accessible without institutional supercomputing infrastructure. Code and data are available at <https://github.com/studiofarzulla/consent-holding-theory>.

The author welcomes feedback, criticism, and collaboration. Correspondence should be directed to [murad@dissensus.ai](mailto:murad@dissensus.ai).

## Declarations

**Conflict of Interest.** The author declares no competing interests.

**Funding.** This research received no external funding.

**Data Availability.** All simulation code and generated data are available at <https://github.com/studiofarzulla/consent-holding-theory>.

**AI Assistance.** Claude (Anthropic) was used as a research collaborator for literature synthesis, computational validation design, LaTeX preparation, and iterative refinement of mathematical arguments. All intellectual claims and errors remain the author’s responsibility.

## References

- Mike Ananny and Kate Crawford. Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability. *new media & society*, 20(3):973–989, 2018.
- Hannah Arendt. *On Revolution*. Viking Press, New York, NY, 1963.
- Kenneth J. Arrow. *Social Choice and Individual Values*. Yale University Press, New Haven, CT, 1951.
- Robert J. Aumann and Roger B. Myerson. Endogenous formation of links between players and of coalitions. *American Economic Review*, 78(2):175–179, 1988.
- John F. Banzhaf III. Weighted voting doesn’t work: A mathematical analysis. *Rutgers Law Review*, 19: 317–343, 1965.
- Sonja Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. 2019. Available at: <https://fairmlbook.org>.

- Wilfred Beckerman and Joanna Pasek. *Justice, Posterity, and the Environment*. Oxford University Press, Oxford, UK, 2001.
- Jenna Bednar. *The Robust Federation: Principles of Design*. Cambridge University Press, Cambridge, UK, 2008.
- Robin Blackburn. *The Overthrow of Colonial Slavery, 1776–1848*. Verso, London, UK, 1988.
- Carles Boix. *Democracy and Redistribution*. Cambridge University Press, Cambridge, UK, 2003.
- Geoffrey Brennan and James M. Buchanan. *The Reason of Rules: Constitutional Political Economy*. Cambridge University Press, Cambridge, UK, 1985.
- Harry Brighouse and Marc Fleurbaey. Democracy and proportionality. *Journal of Political Philosophy*, 18(2):137–155, 2010.
- James M. Buchanan and Gordon Tullock. *The Calculus of Consent: Logical Foundations of Constitutional Democracy*. University of Michigan Press, Ann Arbor, MI, 1962.
- Business Roundtable. Statement on the purpose of a corporation. <https://www.businessroundtable.org/business-roundtable-redefines-the-purpose-of-a-corporation-to-promote-a-n-economy-that-serves-all-americans>, 2019. Accessed November 15, 2025.
- Giovanni Capoccia and R. Daniel Kelemen. The study of critical junctures: Theory, narrative, and counterfactuals. *World Politics*, 59(3):341–369, 2007.
- Erica Chenoweth and Maria J. Stephan. *Why Civil Resistance Works: The Strategic Logic of Nonviolent Conflict*. Columbia University Press, New York, NY, 2011.
- Thomas Christiano. *The Constitution of Equality: Democratic Authority and Its Limits*. Oxford University Press, Oxford, UK, 2008.
- Danielle K. Citron and Saul X. Levmore. The internet, privacy, and you. *Yale Journal of Law & Technology*, 36:405–473, 2014.
- Thomas Clarkson. *The History of the Rise, Progress, and Accomplishment of the Abolition of the African Slave-Trade by the British Parliament*. Longman, Hurst, Rees, and Orme, London, UK, 1808.
- Dimitri Courant and Théo Bourgeron. The french citizens’ convention for climate: Proposals and political outcomes. *French Politics*, 19:321–338, 2021.
- Michael Cox, Gwen Arnold, and Sergio Villamayor Tomás. A review of design principles for community-based natural resource management. *Ecology and Society*, 15(4):38, 2010.
- Robert A. Dahl. *A Preface to Democratic Theory*. University of Chicago Press, Chicago, IL, 1956.
- Robert A. Dahl. *Polyarchy: Participation and Opposition*. Yale University Press, New Haven, CT, 1971.
- Evelyn Douek. Content moderation as systems thinking. *Harvard Law Review*, 136:526–606, 2022.
- Seymour Drescher. *Capitalism and Antislavery: British Mobilization in Comparative Perspective*. Oxford University Press, Oxford, UK, 1987.

- Joshua M. Epstein. *Generative Social Science: Studies in Agent-Based Computational Modeling*. Princeton University Press, Princeton, NJ, 2006.
- Olaudah Equiano. *The Interesting Narrative of the Life of Olaudah Equiano, or Gustavus Vassa, the African*. Author, London, UK, 1789.
- David Estlund. *Democratic Authority: A Philosophical Framework*. Princeton University Press, Princeton, NJ, 2008.
- Christopher J. Fariss. Respect for human rights has improved over time. *Nature News*, 505(7485): 640–642, 2014.
- David M. Farrell, Jane Suiter, and Clodagh Harris. ‘systematizing’ constitutional deliberation: The 2016–18 citizens’ assembly in ireland. *Irish Political Studies*, 34(1):113–123, 2019.
- Murad Farzulla. The axiom of consent: Friction dynamics in multi-agent coordination. *arXiv preprint arXiv:2601.06692*, 2025a. doi: 10.48550/arXiv.2601.06692. Unified friction framework for multi-agent coordination. Code: <https://github.com/studiofarzulla/friction-mar1>.
- Murad Farzulla. ROM: Scale-relative formalism for persistence-conditioned dynamics. *arXiv preprint arXiv:2601.06363*, 2025b. doi: 10.48550/arXiv.2601.06363. Formal foundation for selection-transmission dynamics. Code: <https://github.com/studiofarzulla/consent-rom-empirical1>.
- Larry Fauver and Marie E. Fuerst. Does good corporate governance include employee representation on boards? evidence from german corporate boards. *Journal of Financial Economics*, 99(3):554–569, 2011. doi: 10.1016/j.jfineco.2010.10.016.
- Dan S. Felsenthal and Moshé Machover. *The Measurement of Voting Power: Theory and Practice, Problems and Paradoxes*. Edward Elgar, Cheltenham, UK, 1998.
- Sidney Fine. *Sit-Down: The General Motors Strike of 1936–1937*. University of Michigan Press, Ann Arbor, MI, 1969.
- James S. Fishkin. *When the People Speak: Deliberative Democracy and Public Consultation*. Oxford University Press, Oxford, UK, 2009.
- James S. Fishkin. *Democracy When the People Are Thinking: Revitalizing Our Politics Through Public Deliberation*. Oxford University Press, Oxford, UK, 2018.
- R. Edward Freeman. *Strategic Management: A Stakeholder Approach*. Pitman, Boston, MA, 1984.
- Ehud Friedgut, Gil Kalai, Nathan Keller, and Noam Nisan. A quantitative version of the gibbard-satterthwaite theorem for three alternatives. *SIAM Journal on Computing*, 40(3):684–702, 2011. doi: 10.1137/100791671.
- Milton Friedman. The social responsibility of business is to increase its profits. *New York Times Magazine*, September 13 1970.
- Francis Fukuyama. *Political Order and Political Decay: From the French Revolution to the Present*. Farrar, Straus and Giroux, New York, NY, 2014.

- Michael Gallagher. Electoral systems across the world. JSTOR Library Edition, 2017.
- Stephen M. Gardiner. *A Perfect Moral Storm: The Ethical Tragedy of Climate Change*. Oxford University Press, Oxford, UK, 2011.
- Allan Gibbard. Manipulation of voting schemes: A general result. *Econometrica*, 41(4):587–601, 1973.
- Tarleton Gillespie. *Custodians of the Internet: Platforms, Content Moderation, and the Hidden Decisions That Shape Social Media*. Yale University Press, New Haven, CT, 2018.
- Mark Granovetter. Threshold models of collective behavior. *American Journal of Sociology*, 83(6):1420–1443, 1978. doi: 10.1086/226707.
- Stephan Grimmlikhuijsen, Albert Meijer, Nadine Lieberherr, et al. Legitimacy of algorithmic decision-making: Three value-based strategies. *Public Administration Review*, 82(5):800–812, 2022.
- Jürgen Habermas. *The Theory of Communicative Action, Volume 1: Reason and the Rationalization of Society*. Beacon Press, Boston, MA, 1984. Translated by Thomas McCarthy.
- Jürgen Habermas. *Moral Consciousness and Communicative Action*. MIT Press, Cambridge, MA, 1990. Translated by Christian Lenhardt and Shierry Weber Nicholsen.
- Thomas Hobbes. *Leviathan*. Andrew Crooke, London, UK, 1651.
- Lu Hong and Scott E. Page. Groups of diverse problem solvers can outperform groups of high-ability problem solvers. *Proceedings of the National Academy of Sciences*, 101(46):16385–16389, 2004.
- Intergenerational Foundation. Intergenerational justice and climate change: Comparative policy review. IF, 2024.
- Simon Jäger, Shakked Noy, and Benjamin Schoefer. What does codetermination do? *ILR Review*, 75(4):857–890, 2022.
- Ehud Kalai and Meir Smorodinsky. Other solutions to nash’s bargaining problem. *Econometrica*, 43(3):513–518, 1975. doi: 10.2307/1914280.
- Nathan Keller. A tight quantitative version of arrow’s impossibility theorem. *Journal of the European Mathematical Society*, 14(4):1253–1278, 2012. doi: 10.4171/JEMS/331.
- Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. In *Proceedings of Innovations in Theoretical Computer Science*, 2017.
- Knights of Labor. Preamble and declaration of principles, 1878.
- Christine M. Koggel. Feminist relational theory and global justice. In Thom Brooks, editor, *The Oxford Handbook of Global Justice*, pages 285–302. Oxford University Press, Oxford, UK, 2022.
- Keisuke Koshimizu, Toshiyuki Nakamura, Toshihiro Kamishima, Takanori Murakami, and Takeshi Tezuka. Participatory ai design for fairer algorithmic decision making. In *FAccT 2020: Conference on Fairness, Accountability, and Transparency*, pages 785–795, 2020.
- Hélène Landemore. *Democratic Reason: Politics, Collective Intelligence, and the Rule of the Many*. Princeton University Press, Princeton, NJ, 2013.

- Ethan J. Leib. *Deliberative Democracy in America: A Proposal for a Popular Branch of Government*. Penn State Press, University Park, PA, 2004.
- Arend Lijphart. *Patterns of Democracy: Government Forms and Performance in Thirty-Six Countries*. Yale University Press, New Haven, CT, 2nd edition, 2012.
- John Locke. *Second Treatise of Government*. Hackett, Indianapolis, IN, 1980. Edited by C. B. Macpherson.
- Thomas Lux and Sascha Schiffko. Estimation of agent-based models using sequential monte carlo methods. *Journal of Economic Dynamics and Control*, 91:391–408, 2018.
- Catriona Mackenzie. Three dimensions of autonomy: A relational analysis. In Andrea Veltman and Mark Piper, editors, *Autonomy, Oppression, and Gender*, pages 15–41. Oxford University Press, Oxford, UK, 2014.
- Catriona Mackenzie and Natalie Stoljar, editors. *Relational Autonomy: Feminist Perspectives on Autonomy, Agency, and the Social Self*. Oxford University Press, New York, NY, 2000.
- James Mahoney. *Colonialism and Postcolonial Development: Spanish America in Comparative Perspective*. Princeton University Press, Princeton, NJ, 2010.
- Ewan McGaughey. The codetermination bargains: The history of german corporate and labour law. *Columbia Journal of European Law*, 23(1):135–176, 2016.
- John Stuart Mill. *On Liberty*. John W. Parker and Son, London, UK, 1859.
- John Stuart Mill. *Considerations on Representative Government*. Parker, Son, and Bourn, London, UK, 1861.
- Elchanan Mossel, Krzysztof Oleszkiewicz, and Amartya Sen. Quantitative gibbard-satterthwaite theorem without neutrality. *Combinatorica*, 32(3):305–317, 2012. doi: 10.1007/s00493-012-2713-4.
- John F. Nash. The bargaining problem. *Econometrica*, 18(2):155–162, 1950. doi: 10.2307/1907266.
- Jennifer Nedelsky. Reconceiving autonomy: Sources, thoughts and possibilities. *Yale Journal of Law and Feminism*, 1(1):7–36, 1989.
- Safiya U. Noble. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press, New York, NY, 2018.
- Pippa Norris. Driving democracy: Do power-sharing institutions work? *Journal of Democracy*, 19(4): 14–28, 2009.
- Elinor Ostrom. *Governing the Commons: The Evolution of Institutions for Collective Action*. Cambridge University Press, Cambridge, UK, 1990.
- Vincent Ostrom, Charles M. Tiebout, and Robert Warren. The organization of government in metropolitan areas: A theoretical inquiry. *American Political Science Review*, 55(4):831–842, 1961.
- Fabienne Peter. Political legitimacy. In *Stanford Encyclopedia of Philosophy*. 2010. Available at: <https://plato.stanford.edu/entries/political-legitimacy/>.

- Robert Allen Phillips. *Stakeholder Theory and Organizational Ethics*. Berrett-Koehler, San Francisco, CA, 2003.
- Hanna Fenichel Pitkin. *The Concept of Representation*. University of California Press, Berkeley, CA, 1967.
- Francisco O. Ramirez, Yasemin Soysal, and Suzanne Shanahan. The changing logic of political citizenship: Cross-national acquisition of women's suffrage rights, 1890 to 1990. *American Sociological Review*, 62(5):735–745, 1997.
- John Rawls. *A Theory of Justice*. Harvard University Press, Cambridge, MA, 1971.
- Christian P. Robert and George Casella. *Monte Carlo Statistical Methods*. Springer, New York, NY, 2nd edition, 2004.
- Jean-Jacques Rousseau. *The Social Contract and Other Later Political Writings*. Cambridge University Press, Cambridge, UK, 1997. Edited by Victor Gourevitch.
- Dietrich Rueschemeyer, Evelyn Huber Stephens, and John D. Stephens. *Capitalist Development and Democracy*. University of Chicago Press, Chicago, IL, 1992.
- Mark Allen Satterthwaite. Strategy-proofness and arrow's conditions: Existence and correspondence theorems for voting procedures and social welfare functions. *Journal of Economic Theory*, 10(2): 187–217, 1975.
- Leonard J. Savage. *The Foundations of Statistics*. John Wiley & Sons, New York, NY, 1954.
- Fritz W. Scharpf. *Governing in Europe: Effective and Democratic?* Oxford University Press, Oxford, UK, 1999.
- Vivien A. Schmidt. Democracy and legitimacy in the european union revisited. *Journal of Common Market Studies*, 51(1):131–147, 2013.
- Andrew D. Selbst and Sonja Barocas. The watery garden: Technical and legal opacity in agricultural licensing. *Yale Journal of Law & Technology*, 30:288–350, 2019.
- Amartya Sen. *Development as Freedom*. Oxford University Press, Oxford, UK, 1999.
- Amartya K. Sen. *Collective Choice and Social Welfare: Expanded Edition*. Harvard University Press, Cambridge, MA, 2017.
- Lloyd S. Shapley and Martin Shubik. A method for evaluating the distribution of power in a committee system. *American Political Science Review*, 48(3):787–792, 1954.
- Henry Shue. *Basic Rights: Subsistence, Affluence, and US Foreign Policy*. Oxford University Press, Oxford, UK, 2nd edition, 1993.
- Elizabeth Cady Stanton, Lucretia Mott, and Others. Declaration of sentiments and resolutions. Seneca Falls Convention, July 19–20, 1848, 1848.
- Dawn Langan Teele. How the west was won: Competition, mobilization, and women's enfranchisement in the united states. *Journal of Politics*, 80(2):442–461, 2018.



- Sigurt Vitols. Coordinated market economies and financialization: Germany and japan in comparison. *Competition & Change*, 15(4):288–311, 2011.
- Ari Ezra Waldman and Kristin N. Johnson. The role of decision importance and governance in algorithmic legitimacy. *Northwestern University Law Review*, 117:1–52, 2022.
- Jeremy Waldron. *Law and Disagreement*. Oxford University Press, Oxford, UK, 1999.
- William Wilberforce. Speeches on the abolition of the slave trade. Hansard Parliamentary Debates, 1789–1807.
- Paramita Yadav, Ajay Kumar, and M.N. Raizada. Institutional effects on ecological outcomes of community-based management of fisheries in the amazon: The arapaima case. *PNAS*, 118(29): e2100896118, 2021. doi: 10.1073/pnas.2100896118.