# The Narrow Path Between Domination and Chaos: Consent, Coordination, and Structural Control in the Age of Accelerating Intelligence

Flyxion

January 10, 2026

### Abstract

This paper develops a structural theory of AI governance aimed at eliminating two catastrophic failure modes that arise under accelerating intelligence: domination, in which unilateral power escapes contestation, and chaos, in which multipolar competition erodes coordination through unmanaged externalities. Rather than treating these outcomes as contingent on agent values, incentives, or trust, the paper derives them as structural attractors of systems that permit action without enforceable constraints.

To address this, the paper introduces and formalizes a unified architectural stack—Spherepop OS, Polyxan, and PlenumHub—that reconstitutes governance as an event-sourced, proof-carrying, and externally bounded process. Actions are admitted only through explicit admissibility judgments pairing them with verifiable proofs of compliance against declared specifications, consent boundaries, and externality budgets. Using categorical and sheaf-theoretic semantics, the paper shows that admissibility is compositional, revocation is causal, and coordination is enforced ex ante rather than negotiated ex post. Negative externalities are formalized as failures of global gluing, and are thereby excluded structurally rather than mitigated after the fact.

The core theoretical result is that domination and chaos correspond to objects or trajectories that are not definable within the resulting admissibility space. Domination requires uncertified or non-contestable actions, which cannot exist in a proof-carrying pullback category. Chaos requires locally admissible but globally incoherent action families, which are excluded by sheaf conditions enforced by PlenumHub. These results are supported by rigorous appendices, including a formal grammar, categorical constructions, sheaf-theoretic interpretation, and mechanically checked Lean proofs demonstrating invariant preservation under trajectory extension.

The paper reframes AI alignment as an environmental property of admissibility fields rather than an internal property of agents. It argues that stability under intelligence acceleration is achievable when governance is treated as a conserved structural invariant, and that the so-called "narrow path" between domination and chaos is not a delicate balance but a region made stable by design.

# 1 Introduction

The dominant failures of contemporary digital platforms are often framed as problems of moderation, misinformation, or malicious actors. In the context of artificial intelligence, these concerns escalate into familiar apocalyptic narratives: either a singular system becomes uncontrollable and subjugates humanity, or a fragmented landscape of competing systems erodes social order beyond repair. While these narratives differ in tone, they share a common limitation. They treat civilizational failure as the result of bad outcomes rather than bad structures.

This paper advances a different diagnosis. The central problem is not the intelligence of systems, nor even their objectives, but the absence of architectural constraints capable of preserving human agency under conditions of accelerating optimization. When systems grow faster than the institutions, epistemic tools, and consent mechanisms designed to govern them, failure becomes structural rather than accidental. In such regimes, outcomes are no longer meaningfully chosen; they are selected by dynamics that operate beyond human steering capacity.

Recent work by Nora Ammann has articulated this situation as a "narrow path" between two civilizational attractors: domination and chaos. Domination refers to scenarios in which power becomes so concentrated, whether in a single system or a small coalition, that meaningful contestation is no longer possible. Chaos refers to multipolar regimes in which excessive competition destroys coordination, forcing rational actors into arms races that erode shared values, externalize risk, and collapse trust. The force of Ammann's contribution lies not in predicting specific disasters, but in identifying the structural corridor within which civilization must remain if it is to survive the transition to machine-mediated intelligence.

This essay reconstructs that framework from more general theoretical foundations concerning consent, engagement, optimization, and semantic drift. Rather than treating domination and chaos as contingent outcomes to be avoided through better intentions or more cautious actors, they are derived here as necessary consequences of unconstrained optimization operating on human attention, identity, and meaning. From this perspective, the narrow path is not a metaphor but a feasible region in sociotechnical design space.

The second aim of the paper is constructive. Drawing on my own work in event-sourced computation, semantic operating systems, and proof-carrying coordination infrastructures, I argue that the missing ingredients are neither technical capability nor user demand, but explicit architectural layers for consent, constraint, and verification. These layers already exist in other domains of computing and governance. Their absence in social and AI-driven platforms is not an inevitability, but a choice shaped by incentive structures.

The paper proceeds as follows. The next section establishes a conceptual framework by defining consent, engagement, optimization objectives, and semantic drift with sufficient precision to support formal analysis. Subsequent sections examine empirical evidence of systemic failure, analyze the technical mechanisms that generate these outcomes, and situate generative AI as an accelerant rather than a root cause. Later sections compare alternative architectures, develop concrete design proposals, address objections, and explore regulatory and philosophical implications. The conclusion reframes

the problem as a general principle of alignment: optimization without constraint inevitably violates autonomy.

## 2 Conceptual Framework

Any rigorous analysis of platform behavior and AI-mediated systems must begin with precise definitions of the terms that are most frequently invoked and least clearly specified. Concepts such as consent, engagement, and optimization are often treated as self-evident, yet their ambiguity allows fundamentally different phenomena to be conflated under a single label. This section establishes a conceptual framework that distinguishes these notions and clarifies the mechanisms by which misalignment emerges.

Consent, in its strongest and most widely accepted ethical formulations, is not a one-time act of acquiescence but an ongoing, informed, and revocable authorization. In medical ethics and human-subjects research, consent requires that participants understand the nature of what they are agreeing to, retain the ability to withdraw without penalty, and are not coerced by asymmetries of power or information. By contrast, consent in digital platforms is typically reduced to acceptance of generalized terms of service that are neither specific to actual experiences nor meaningfully revisitable once granted.

This paper adopts a definition of consent as the sustained ability of a user to declare boundaries over categories of experience, to have those declarations respected by default, and to revise them over time with predictable effect. Under this definition, consent is violated not only when harmful content is shown, but when users are systematically deprived of the means to specify what they do not wish to encounter. The absence of such mechanisms renders subsequent claims of voluntary engagement ethically hollow.

Engagement, similarly, must be disambiguated. Meaningful engagement refers to deliberate interaction that users value upon reflection, even if it requires effort or discomfort. Extracted engagement refers to interaction induced by compulsion, novelty, or emotional manipulation that users later regret or would prefer to avoid. Both forms may be behaviorally indistinguishable at the level of clicks, dwell time, or shares, yet they differ fundamentally in their relationship to autonomy and welfare.

Current platform metrics overwhelmingly conflate these categories. Optimization systems treat all engagement as positive signal, regardless of whether it corresponds to satisfaction, distress, or mere reflex. This conflation is not a bug but a structural feature of inference systems that rely exclusively on behavioral traces without access to explicit user intent.

Optimization objectives provide the formal substrate on which these dynamics operate. Platforms do not optimize for abstract engagement but for specific proxy metrics such as click-through rates, session duration, frequency of return, and advertising yield. These metrics form a hierarchy in which revenue objectives dominate user welfare by design, even when public narratives emphasize personalization or connection. The resulting objective functions are mathematically coherent but ethically indifferent, favoring content that reliably captures attention regardless of long-term effects.

Semantic drift describes the process by which content categories lose their descriptive fidelity

under sustained optimization pressure. Formally, let C denote a set of content labels and F denote the functional effects those contents have on users, such as arousal, distress, or compulsion. Under engagement maximization, the mapping from C to F diverges over time as producers learn to exploit classification boundaries. Content nominally labeled as entertainment, lifestyle, or humor increasingly converges toward the same underlying affective triggers, even as surface features vary.

This divergence is accelerated by machine learning systems that optimize for outcomes rather than meanings. As generative models enable infinite variation at negligible cost, semantic drift becomes a dominant force, eroding users' ability to predict or control their experience based on categories alone. Without corrective constraints, the informational environment collapses toward lowest-common-denominator attractors that maximize engagement while hollowing out semantic distinctions.

Together, these definitions establish the conditions under which domination and chaos arise. When consent is reduced to acquiescence, engagement is conflated with welfare, optimization is unconstrained, and semantics are allowed to drift, systems naturally evolve toward regimes that either centralize power irreversibly or fragment coordination beyond repair. The narrow path consists in restoring explicit constraints at each of these layers.

## 3 Evidence and Patterns

The structural claims advanced in the previous sections are not speculative. They are supported by a growing body of empirical evidence drawn from user surveys, platform transparency reports, academic studies, and longitudinal analyses of platform behavior. When examined collectively, these sources reveal a consistent pattern: users experience persistent exposure to unwanted content, blocking and feedback mechanisms fail to generalize, and content ecosystems drift toward increasingly extreme or intrusive forms under engagement optimization.

Large-scale surveys conducted across North America and Europe between 2018 and 2024 consistently report that between 55% and 72% of users encounter content they explicitly find unwanted or distressing on major social platforms at least weekly, with 30% to 45% reporting daily exposure. These figures are remarkably stable across age groups, though adolescents and young adults report higher exposure to sexualized and violent material, while older users report higher exposure to scams and political manipulation. Importantly, a majority of respondents in these studies report that they have actively attempted to reduce such exposure through available tools, including blocking, hiding, and content reporting.

Platform transparency reports corroborate these findings indirectly. Meta's quarterly Community Standards Enforcement Reports show a year-over-year increase in content removals for spam, sexual solicitation, and deceptive practices, despite simultaneous claims of improved moderation. The absolute volume of removed content increases even as enforcement efficiency improves, indicating that content generation and distribution are growing faster than moderation capacity. Similar trends are observable in TikTok's Integrity Reports and Google's Transparency Reports for YouTube.

Qualitative case analyses further illuminate the failure modes. Sexualized clickbait content, for

example, frequently appears in forms that remain formally compliant with platform policies while exploiting visual framing, suggestive motion, or symbolic proxies. Content creators systematically test boundary conditions, producing variants that evade automated classifiers while retaining high engagement performance. Engagement metrics for such content routinely exceed platform averages by factors of two to five, particularly in short-form video formats.

User blocking behavior provides a second line of evidence. Studies analyzing anonymized interaction logs indicate that users often block dozens or even hundreds of accounts over time, yet report little perceived improvement in feed quality. This is explained by the fact that blocking operates at the level of individual accounts rather than content categories or generative patterns. When blocked content is algorithmically replaced by near-identical material from different sources, blocking functions as a local perturbation rather than a structural constraint.

Complaint and feedback systems exhibit similar limitations. Reports are processed on a per-item basis, with no durable effect on downstream recommendation policies for the reporting user. As a result, users engage in repeated micro-interventions that fail to accumulate into meaningful preference signals. Over time, this leads to consent fatigue, in which users disengage from feedback mechanisms altogether while remaining subject to the same exposure patterns.

Temporal analysis shows that these dynamics have intensified rather than stabilized. The introduction of large-scale generative models between 2021 and 2024 corresponds with a marked increase in low-cost, high-volume content production. Detection evasion techniques, such as slight visual variation or paraphrasing, defeat duplicate detection systems, allowing spam and clickbait to scale exponentially while moderation resources grow linearly.

Cross-platform comparison strengthens the argument that these outcomes are architectural rather than incidental. Platforms that rely primarily on algorithmic feeds optimized for engagement, such as TikTok, Instagram Reels, and Facebook Feed, exhibit higher reported rates of unwanted content exposure than systems organized around explicit subscription or search paradigms. By contrast, RSS-based systems, email clients with rule-based filters, and academic databases with boolean query operators allow users to express durable negative preferences that generalize across sources.

Taken together, these patterns support a strong empirical conclusion: current platform architectures systematically fail to respect user consent at scale, not because of insufficient moderation effort, but because they lack mechanisms for expressing and enforcing category-level boundaries under optimization pressure.

## 4 Technical Architecture Analysis

To understand why these failures are persistent, it is necessary to examine the technical architecture of contemporary recommendation systems and the formal structure of their optimization objectives. While implementations vary, most large platforms employ a combination of collaborative filtering, content-based recommendation, and deep learning embeddings trained on historical interaction data. These systems are unified by a common mathematical structure.

Let $U$ denote the set of users, $C$ the set of content items, and $t$ discrete time steps. Let $E(c, u, t)$

denote a scalar engagement signal measuring the interaction of user $u$ with content $c$ at time $t$, where engagement may aggregate clicks, dwell time, reactions, or shares. Let $P(c \mid u, t)$ denote the probability that content $c$ is shown to user $u$ at time $t$.

The canonical optimization objective can be written as

$$\max_{P} \sum_{u \in U} \sum_{c \in C} \sum_{t} E(c, u, t) \cdot P(c \mid u, t),$$

subject to policy constraints that exclude explicitly prohibited content classes.

This formulation reveals the core problem. Engagement signals are treated as sufficient statistics for user preference, and optimization operates exclusively through inferred behavior rather than declared intent. The system has no access to the semantic meaning of engagement, only its magnitude.

We can formalize the resulting feedback loop as a dynamical system. Let $\theta_t$ denote the parameters of the recommendation model at time $t$. Then

$$\theta_{t+1} = \theta_t + \eta \nabla_\theta \mathbb{E}_{u,c}[E(c, u, t)],$$

where $\eta$ is a learning rate. Content producers adapt symmetrically by generating content that maximizes expected engagement under the current model. This coupled optimization creates an amplification cycle in which both sides converge toward engagement-maximizing attractors.

Proposition 1: Under engagement maximization without explicit user-declared constraints, content distributions converge toward stimuli that maximize short-term affective response rather than reflective user preference.

Proof sketch: Let $F(c)$ denote the affective impact of content $c$, and let $R(c)$ denote reflective user valuation. Empirical evidence shows that $E(c, u, t)$ correlates more strongly with $F(c)$ than with $R(c)$ when measurement is immediate. Since optimization maximizes $E$, gradients favor increases in $F$ regardless of $R$. Over repeated iterations, content with high $F$ but low $R$ dominates the distribution. The absence of corrective constraints allows this divergence to persist.

Information asymmetry further exacerbates the problem. Users communicate preferences through low-bandwidth signals such as clicks or hides, while platforms infer high-dimensional behavioral models incorporating thousands of latent features. The mutual information between user intent and observable behavior is strictly bounded, while platform inference capacity is not. This asymmetry makes it formally impossible to reliably distinguish curiosity from interest or disgust from attraction using behavior alone.

Proposition 2: No behavior-only inference system can reliably recover negative preferences at scale.

Proof: Let $X$ be the latent variable representing user intent and $Y$ the observed behavior. If two distinct intents $x_1 \neq x_2$ induce overlapping distributions over $Y$, then no classifier can distinguish them with zero error. Empirical studies demonstrate such overlap for curiosity, outrage, and aversion. Therefore, negative preference inference is provably lossy.

This impossibility result implies that any architecture relying solely on inferred engagement is

structurally incapable of respecting consent. The failure is not due to insufficient model capacity, but to missing input channels.

In the next section, this analysis is extended to generative AI systems, which dramatically accelerate these dynamics by reducing content production costs and eliminating identity friction.

## 5  Generative AI, Identity Friction, and Accelerated Failure Modes

The introduction of large-scale generative models fundamentally alters the dynamics analyzed in the previous section, not by changing the structure of optimization, but by dramatically reducing the cost of exploiting it. Generative adversarial networks, diffusion models, and large language models enable the production of arbitrarily many content variants at negligible marginal cost. This removes a key historical constraint on spam, clickbait, and manipulative content: the labor required to produce variation.

Let $G(\cdot)$ denote a generative model parameterized to produce content $c$ conditioned on prompts or latent variables. The expected cost of producing $n$ variants scales sublinearly with $n$, while the probability that at least one variant evades automated detection increases monotonically with $n$. As a result, content producers face a convex payoff curve in which small increases in generation volume yield disproportionately large increases in distribution success.

Empirical estimates place the cost of generating one thousand image or video variants at orders of magnitude below the cost of equivalent human-produced content. At the same time, moderation systems rely on pattern recognition and duplicate detection that are brittle under small perturbations. Slight changes in framing, color, cropping, or phrasing are sufficient to evade filters without altering affective impact.

This asymmetry produces an exponential growth dynamic. Let $V_t$ denote the volume of generated content at time $t$, and let $M_t$ denote moderation capacity. If $V_t = V_0 e^{\alpha t}$ while $M_t = M_0 + \beta t$, then for any $\alpha > 0$, there exists a finite $t^*$ beyond which moderation is structurally overwhelmed. This is not a failure of moderation quality but a mismatch of growth rates.

## 6  Economic Incentives and the Profitability of Abuse

The economic logic underpinning these dynamics can be formalized directly. Consider an operator deploying AI-generated content across a platform. Let $R$ denote expected revenue per engagement, $E$ expected engagement volume, $C_g$ generation cost, and $C_a$ expected account loss cost due to enforcement. The expected value of an operation is given by

$$\mathbb{E}[V] = R \cdot E - (C_g + C_a).$$

Under generative AI, $C_g$ approaches zero, and $C_a$ is amortized across large numbers of disposable accounts. As long as $R \cdot E$ remains positive, abuse becomes economically rational. Engagement optimization by platforms increases $E$ precisely for the content types most likely to exploit affective

triggers, further increasing expected value.

This explains the observed proliferation of AI-generated spam, scams, and identity farms. Removing friction from content creation does not merely increase volume; it shifts entire classes of behavior from marginal to dominant strategies. Content that would previously have been unprofitable due to labor costs becomes ubiquitous.

# 7   Identity Persistence and Accountability Collapse

Identity functions as a deterrence mechanism in social systems by linking actions to long-term reputational consequences. When identity is persistent, harmful behavior incurs future costs. When identity is disposable, deterrence collapses.

Generative AI enables the rapid creation of synthetic personas at scale. These personas can be discarded and replaced faster than enforcement systems can react. The expected penalty for abusive behavior approaches zero as identity half-life decreases. This phenomenon is observable across platforms in the form of bot swarms, cloned influencer accounts, and rotating scam identities.

Formally, let $I$ denote identity persistence measured as expected lifespan. Deterrence requires that the expected cost of sanctions exceeds expected gain:

$$\mathbb{E}[\text{Sanction}] \cdot I > \mathbb{E}[\text{Gain}].$$

As $I \to 0$, this inequality cannot hold. No amount of content moderation can restore deterrence in the absence of persistent identity binding.

Systems that retain accountability, such as academic publishing, financial institutions, and professional licensing regimes, all rely on strong identity persistence and traceability. Social platforms, by contrast, intentionally weaken identity constraints to maximize participation and growth, thereby eliminating the very mechanism that would discourage abuse at scale.

# 8   From Engagement Pathologies to Civilizational Attractors

The preceding analyses establish that engagement-optimized systems without explicit consent constraints produce predictable pathologies at the level of content distribution, identity, and economic behavior. The final step is to show how these pathologies scale into the civilizational failure modes described as domination and chaos.

Domination arises when optimization and acceleration concentrate power faster than contestation mechanisms can respond. In AI-mediated systems, this occurs when control over infrastructure, models, or coordination layers becomes irreversible due to speed, opacity, or scale. Once a system or coalition can act beyond the temporal reach of oversight, human intervention becomes symbolic rather than causal.

Chaos arises when the same acceleration occurs in a multipolar environment without shared constraints. Actors race to deploy increasingly powerful systems defensively, anticipating domination

by others. Coordination collapses not because actors disagree on values, but because unilateral restraint becomes irrational. Externalities such as environmental damage, security risk, and epistemic degradation are ignored in favor of survival.

Proposition 3: Under accelerating optimization without enforceable global constraints, sociotechnical systems converge toward either domination or chaos.

Proof sketch: Let $S$ denote system capability growth rate, and let $C$ denote coordination capacity. If $S > C$ and power is centralized, domination results. If $S > C$ and power is distributed, chaos results. The absence of constraint ensures $S$ grows unchecked, while $C$ grows slowly due to institutional inertia. Therefore convergence toward one of the two attractors is inevitable.

The narrow path corresponds to maintaining $C \geq S$ through structural interventions that scale coordination and oversight alongside capability.

## 9 The Narrow Path as a Structural Invariant

The narrow path is not a compromise between extremes but a region of parameter space in which human agency remains viable. Its defining feature is not moderation, but constraint. Consent mechanisms, proof-carrying systems, and defense-favored architectures function as invariants that prevent drift toward irreversible states.

This framing reveals why appeals to ethics, trust, or goodwill are insufficient. Without enforceable structure, even aligned actors are swept along by optimization dynamics. Alignment must therefore be understood as a property of architectures, not intentions.

## 10 Design Alternatives and Constraint-Based Architectures

The failures described thus far are not inherent to computation or intelligence. They arise from specific architectural omissions. Systems in other domains demonstrate that explicit constraint layers can coexist with powerful optimization.

Declarative filtering systems such as email rules, subscription-based feeds, and query-driven databases allow users to specify both positive and negative preferences. These preferences are treated as hard constraints rather than soft signals. Blocking generalizes across sources because it operates at the level of categories, not instances.

A consent-respecting architecture introduces an explicit constraint function $U(c) \in \{0, 1\}$ representing user-declared admissibility. The optimization objective becomes

$$\max_P \sum E(c, u, t) \cdot P(c \mid u, t) \quad \text{subject to} \quad U(c) = 1.$$

This modification is trivial computationally yet transformative structurally. It restores agency by making boundaries durable, predictable, and enforceable.

Such architectures also enable auditability. Violations can be detected mechanically, and compliance can be verified externally. This is a prerequisite for meaningful governance.

## 11    Consent Architecture and the Limits of Implied Agreement

The failures of contemporary platforms cannot be fully understood without examining the architecture of consent itself. Consent, as implemented in most digital systems, is treated as a one-time legal artifact rather than an ongoing operational constraint. Users are asked to accept terms of service that describe broad categories of data use and content exposure, yet these documents neither enumerate the specific experiential consequences of participation nor provide mechanisms for granular withdrawal.

In contrast, established ethical frameworks in medicine and human-subjects research impose far more stringent requirements. The Belmont Report, for example, defines respect for persons as requiring informed consent, comprehension, and voluntariness. Consent must be specific to the intervention, revisable over time, and withdrawable without undue burden. These standards exist precisely because asymmetries of knowledge and power render blanket consent ethically insufficient.

Digital platforms systematically violate these principles. Users cannot meaningfully anticipate the content they will be exposed to, nor can they revoke consent for entire classes of experience once participation has begun. The gap between what users nominally agree to and what they actually encounter constitutes a structural consent failure rather than a misunderstanding.

Withdrawal of consent is particularly instructive. In ethical research contexts, withdrawal is immediate and complete. In platform contexts, withdrawal is fragmented into micro-actions such as blocking individual accounts or hiding specific posts. These actions impose cognitive and emotional costs on users while failing to generalize. As a result, consent withdrawal is both burdensome and ineffective, leading to what may be termed consent fatigue.

Consent fatigue undermines agency by transforming boundary-setting into an exhausting, never-ending task. Each unwanted exposure demands a new response, yet none of these responses accumulate into a stable preference representation. Over time, users disengage from the act of consent itself while remaining subject to the same optimization dynamics.

Power asymmetry exacerbates this problem. Platforms possess detailed knowledge of content effects, engagement patterns, and system behavior. Users, by contrast, experience only their local feed and lack visibility into distribution logic. Asking users to consent under such conditions is analogous to asking patients to consent to treatment without disclosing risks. The asymmetry invalidates the consent.

## 12    Generative AI and the Collapse of Meaningful Agency

The introduction of generative AI compounds consent failures by accelerating semantic drift and identity fragmentation. When content can be generated infinitely and identities can be instantiated at will, users lose not only control over what they see but also the ability to reason about why they are seeing it.

Generative models decouple content appearance from human intention. Previously, encountering objectionable content implied the existence of a producer who could be blocked, reported, or

sanctioned. Under generative regimes, content becomes a fluid surface phenomenon, detached from durable agents. Blocking one instance does nothing to prevent recurrence because there is no stable source to constrain.

This decoupling erodes accountability at every level. Users cannot target their objections effectively, platforms cannot rely on reputational deterrence, and regulators cannot trace responsibility. Agency dissolves into statistical inference, and consent becomes inoperative.

## 13   Assurance, Proof, and Justified Confidence

Restoring agency under these conditions requires more than improved moderation. It requires epistemic infrastructure capable of supporting justified confidence in system behavior. Nora Ammann emphasizes high-assurance systems in which AI outputs are accompanied by verifiable proofs that specified properties have been satisfied. This approach shifts governance from trust to verification.

Formal assurance is not a demand for perfection. It is a recognition that unchecked complexity defeats human oversight. Proof-carrying artifacts compress reasoning into forms that are easier to check than to generate. When verification costs scale sublinearly relative to system capability, oversight remains viable even as systems accelerate.

In software engineering, this principle is well established. Formally verified kernels, memory-safe languages, and proof-carrying code allow systems to enforce invariants that would otherwise be violated by optimization pressure. The same logic applies to sociotechnical systems. Consent constraints, identity bindings, and coordination rules must be enforced mechanically rather than normatively.

Let $S$ denote a system that produces outputs $o$ along with proofs $\pi$ such that $\pi \vdash P(o)$, where $P$ is a formally specified property. Governance consists not in evaluating $o$ directly, but in verifying $\pi$. This indirection is essential for scalability.

## 14   Defense-Favored Structures and Civilizational Stability

A system is defense-favored if small investments in protection suffice to neutralize large offensive efforts. Defense-favoring is not militarization; it is asymmetry design. In cybersecurity, memory safety dramatically reduces exploitability. In content systems, hard consent constraints dramatically reduce abuse profitability.

Defense-favored architectures stabilize civilization during periods of rapid technological transition by raising the cost of exploitation relative to the cost of compliance. They do not eliminate conflict, but they prevent conflict from scaling uncontrollably.

Critically, defense-favoring must apply at the architectural level. Relying on human vigilance or reactive enforcement produces attacker advantages. Embedding constraints into system design reverses this asymmetry.

# 15  From Structural Control to Human Agency

The cumulative effect of these interventions is not total control but restored agency. Humans remain in the loop not by supervising every action, but by specifying the invariants within which action occurs. Consent, coordination, and assurance become structural properties rather than afterthoughts.

If these harnesses are not built before AI systems acquire autonomous research and deployment capabilities, human intervention will persist only as ritual. The opportunity to shape the trajectory of intelligence is therefore temporally bounded.

# 16  Objections and Responses

Any proposal to introduce explicit consent constraints and structural limits on optimization encounters a predictable set of objections. These objections recur across industry discourse, policy debates, and academic commentary. Addressing them rigorously is essential, not because they are novel, but because they reveal the underlying assumptions that sustain the current architecture.

A common objection holds that users do not know what they want, and that algorithmic inference is therefore necessary to surface relevant content. This argument conflates discovery with violation. While users may not be able to articulate every positive preference in advance, they are demonstrably capable of identifying categories of content they wish to avoid. Survey data consistently show high confidence and low regret for negative preferences, particularly regarding sexualized material, scams, violence, and repetitive spam. The inability to predict all future interests does not imply consent to unbounded intrusion. Discovery mechanisms can coexist with hard exclusion constraints without contradiction.

A second objection claims that introducing consent constraints would reduce engagement and threaten platform revenue. This argument treats short-term engagement volume as synonymous with long-term platform viability. However, longitudinal studies of user satisfaction indicate that trust erosion and content fatigue correlate strongly with churn, reduced sharing, and disengagement over time. Engagement extracted through compulsion undermines the very user base on which monetization depends. Moreover, platforms that offer greater user control, such as email clients and subscription-based services, demonstrate stable long-term usage despite lower raw engagement metrics. The relevant comparison is not between maximal engagement and minimal engagement, but between extractive growth and sustainable participation.

A third objection asserts that content classification is too difficult to support category-level controls. This claim is empirically false. Platforms already classify content extensively for advertising targeting, age restrictions, and policy enforcement. The same taxonomies used to sell advertisements can be used to enforce user-declared exclusions. The difficulty lies not in classification, but in the willingness to allow those classifications to constrain distribution.

Concerns about filter bubbles are similarly misplaced. Filter bubbles arise when users are confined to narrow informational regimes without awareness or control. Explicit exclusions chosen by users are categorically different. Choosing not to see sexualized clickbait does not entail ideologi-

cal isolation. On the contrary, agency over boundaries enables more intentional exploration within remaining spaces.

Finally, free speech objections misunderstand the nature of the proposal. Restricting what must be shown to a given user does not restrict what may be said. Content rating systems in broadcasting and parental controls in software demonstrate that choice-enhancing constraints expand freedom rather than diminish it. The absence of consent mechanisms compels exposure, which is itself a form of coercion.

# 17 Policy and Governance Implications

The structural failures described in this paper expose significant gaps between existing legal frameworks and platform practices. Data protection regulations such as the General Data Protection Regulation emphasize informed consent, purpose limitation, and the right to withdraw consent. However, these principles are rarely applied to content exposure itself, despite its profound impact on autonomy and well-being.

Current law largely treats content exposure as an incidental consequence of platform use rather than a regulated activity. This framing is increasingly untenable. When platforms algorithmically curate experiential environments, they function as information stewards with duties analogous to those of medical researchers, broadcasters, or financial intermediaries.

Regulatory approaches can be broadly categorized into disclosure-based, process-based, outcome-based, and audit-based models. Disclosure alone is insufficient, as users cannot act on information without corresponding control mechanisms. Outcome-based regulation risks perverse incentives and censorship pressures. Process-based requirements mandating the provision of durable consent controls offer a more promising path. Audit-based oversight, supported by formal verification and external review, can ensure compliance without dictating content outcomes.

Internationally, emerging frameworks such as the European Union's Digital Services Act gesture toward systemic risk management but stop short of mandating user-level consent architectures. Without such mandates, platforms remain free to optimize around engagement while treating consent as a legal formality. Aligning regulation with architectural requirements is therefore a central governance challenge.

# 18 Philosophical Foundations: Autonomy, Attention, and Information Ethics

The deeper significance of consent architecture lies in its connection to autonomy. Philosophical accounts of autonomy emphasize the capacity for self-governance under conditions of adequate information and freedom from coercion. When algorithmic systems shape the environment of choice itself, autonomy cannot be reduced to isolated decisions within that environment.

Attention is a finite cognitive resource and a prerequisite for agency. When attention is systematically captured, fragmented, or manipulated without consent, individuals are deprived of the

conditions necessary for reflective judgment. The ethics of attention thus form a foundational layer of digital ethics, intersecting with concerns about addiction, persuasion, and manipulation.

Information ethics further clarifies the obligations of platforms as stewards of the infosphere. To curate information environments is to exercise power over meaning, salience, and perception. Such power carries duties of care that cannot be discharged through disclaimers alone. Control over semantic exposure is a core component of informational self-determination.

Structural coercion arises when exit costs are high and alternatives are scarce. Network effects bind users to platforms despite dissatisfaction, undermining the voluntariness of continued participation. In such contexts, the absence of consent mechanisms constitutes a form of compelled exposure incompatible with autonomy.

## 19  From Extraction to Alignment

The failures of contemporary platforms and the risks posed by accelerating AI systems share a common root: optimization without constraint. Engagement maximization, identity fluidity, and semantic drift are not isolated defects but mutually reinforcing consequences of architectures that prioritize growth over agency.

The analysis presented here demonstrates that domination and chaos are not speculative futures but structural attractors toward which unconstrained systems naturally evolve. Domination emerges when optimization concentrates power beyond contestation. Chaos emerges when optimization fragments coordination beyond repair. The narrow path between them is defined by the preservation of consent, accountability, and enforceable invariants.

The solution is not moderation alone, nor better intentions, nor slower progress. It is the introduction of explicit architectural layers that allow users to declare boundaries, systems to prove compliance, and societies to scale coordination alongside capability. These mechanisms are technically feasible, economically justifiable, and ethically required.

The broader implication extends beyond social platforms. Any system that optimizes behavior at scale without incorporating explicit constraints on consent risks violating autonomy by design. As artificial intelligence permeates domains such as hiring, lending, healthcare, and governance, the lessons of content systems become universal.

Alignment, properly understood, is not a matter of teaching machines what to want. It is a matter of ensuring that no system, however intelligent, can override the conditions under which human choice remains meaningful. The narrow path is real, but it is navigable—if its constraints are treated not as obstacles, but as the very structure of freedom.

## 20  Spherepop OS as an Event-Sourced Consent Substrate

The preceding analysis establishes that consent, coordination, and assurance must be enforced at the architectural level rather than inferred post hoc from behavior. Spherepop OS provides a formal substrate for such enforcement by replacing state-based interaction models with event-sourced,

irreversible semantic construction.

Spherepop OS is grounded in the principle that meaning, identity, and authority emerge from histories of irreversible events rather than mutable state snapshots. Formally, let $E = \{e_1, e_2, ...\}$ denote a partially ordered set of events, where each event represents a committed semantic action. System state is not primitive but derived as a fold over event history. This construction ensures that all meaningful change is auditable, ordered, and attributable.

Consent within Spherepop OS is not represented as a mutable preference vector but as a class of boundary-setting events. A consent declaration is an irreversible commitment that constrains all future admissible events within a given semantic scope. Let $C \subset E$ denote consent events, and let $\mathscr{A}(e)$ denote the set of future events admissible after event $e$. Then for any consent event $c \in C$, admissibility is restricted such that

$$\forall e' \succ c, \quad e' \in \mathscr{A}(c) \iff e' \notin \mathscr{F}_c,$$

where $\mathscr{F}_c$ is the forbidden semantic region declared by the consent event.

This formulation resolves the consent problem identified earlier. Boundaries are explicit, durable, and generalize across instances. Withdrawal of consent is itself an event, restoring admissibility within a new scope without erasing history. Consent fatigue is avoided because boundaries persist until explicitly revised, rather than requiring repeated micro-actions.

Spherepop OS further enforces semantic stability by disallowing silent reinterpretation of categories. Semantic drift is prevented by requiring that any change in category meaning be enacted through explicit transformation events that are themselves subject to review and contestation. This restores predictability to the semantic environment and prevents optimization from exploiting category ambiguity.

## 21    Polyxan as a Proof-Carrying Coordination Layer

While Spherepop OS provides event-level constraint and auditability, large-scale coordination requires additional guarantees concerning correctness, safety, and compliance. Polyxan addresses this requirement by formalizing coordination as proof-carrying interaction between agents and systems.

Polyxan is defined as a typed specification and proof system in which actions are only admissible if accompanied by machine-checkable evidence that they satisfy declared invariants. Let $a$ denote an action proposal, $\pi$ a proof object, and $\Phi$ a specification predicate. An action is admissible if and only if

$$\pi \vdash \Phi(a).$$

Crucially, Polyxan separates the generation of proofs from their verification. Verification is designed to be computationally cheaper than proof construction, ensuring that oversight scales favorably relative to action complexity. This asymmetry is essential for maintaining control under accelerating capability.

In the context of AI-mediated systems, Polyxan enables high-assurance delegation. Autonomous agents may propose actions or artifacts, such as code deployments or content transformations, but

these proposals are inert unless accompanied by proofs that they respect consent constraints, safety properties, and coordination rules encoded in $\Phi$.

Polyxan thus instantiates Ammann's notion of justified confidence. Confidence is justified not by trust in agent intent, but by verifiable evidence that constraints have been respected. Governance shifts from discretionary approval to invariant enforcement.

## 22  PlenumHub and Multipolar Coordination Without Chaos

The challenge of avoiding chaos in a multipolar environment is fundamentally a coordination problem under adversarial incentives. PlenumHub is designed as a shared coordination fabric that enables positive-sum interaction while managing externalities explicitly.

PlenumHub operates as a plenum in the literal sense: a shared space in which commitments, constraints, and proofs are published and verified collectively. Agents interacting through PlenumHub do not rely on bilateral trust but on shared verification infrastructure. Let $H$ denote the hub, and let $\mathscr{P}_i$ denote the policy set of agent $i$. Interactions are mediated such that

$$\forall i, j, \quad \text{interaction}_{i,j} \in H \iff \exists \pi \text{ such that } \pi \vdash (\mathscr{P}_i \wedge \mathscr{P}_j).$$

This formulation ensures that coordination does not collapse under defection pressure. Agents cannot unilaterally externalize risk without violating publicly verifiable constraints. Externalities are managed by encoding them into shared invariants rather than relying on after-the-fact enforcement.

PlenumHub further supports defense-favored equilibria by making exploitation visible and provable. When violations are detectable at low cost, attackers lose their asymmetrical advantage. This stabilizes cooperation even in the presence of competitive incentives.

## 23  Integrating the Stack: Structural Alignment by Construction

Taken together, Spherepop OS, Polyxan, and PlenumHub constitute a vertically integrated alignment stack. Spherepop OS enforces irreversible semantic order and consent at the event level. Polyxan enforces correctness and safety through proof-carrying actions. PlenumHub enforces coordination and externality management across agents.

Formally, the stack can be viewed as enforcing a hierarchy of invariants. Let $I_1$ denote semantic and consent invariants, $I_2$ denote action-level correctness invariants, and $I_3$ denote coordination invariants. An action is globally admissible if and only if

$$\pi_1 \vdash I_1 \ \wedge \ \pi_2 \vdash I_2 \ \wedge \ \pi_3 \vdash I_3.$$

This construction directly addresses the failure modes identified earlier. Domination is prevented because no agent can bypass invariants unilaterally. Chaos is prevented because coordination is enforced structurally rather than normatively. Optimization remains possible, but only within a constrained and auditable space.

## 24  Deriving the Narrow Path as a Formal Consequence

With this stack in place, the narrow path ceases to be a fragile balance and becomes a formal consequence of system design. Capability growth no longer automatically implies loss of control because steering capacity scales with verification rather than human attention. Competition no longer implies chaos because externalities are internalized through shared constraints.

Proposition 4: In a system where all actions are event-sourced, proof-carrying, and hub-coordinated, neither domination nor chaos is a stable attractor.

Proof sketch: Domination requires unilateral action beyond contestation, which is excluded by invariant enforcement. Chaos requires incentive structures that reward defection, which is excluded by shared verification and externality encoding. Therefore, both attractors are eliminated by construction.

This result does not assume benevolence, alignment of values, or absence of conflict. It assumes only that invariants are enforced mechanically and verification remains cheaper than exploitation.

## 25  Event-Sourced Semantics and the Failure of State-Based Governance

The dominant computational paradigm underlying contemporary digital platforms is state-based. Systems are modeled as collections of mutable variables updated over time, with meaning inferred from the current configuration rather than from the path by which it was reached. While this paradigm is efficient for optimization, it is fundamentally hostile to accountability, consent, and governance under acceleration.

In a state-based system, history is optional. Past actions may be logged, summarized, or discarded without affecting current behavior. This makes it impossible to enforce invariants that depend on irreversible commitments. Consent, when represented as a mutable flag, can be silently overridden or reinterpreted without detectable violation. Identity, when represented as a pointer, can be reassigned or duplicated without consequence.

Spherepop OS rejects this paradigm by treating events, rather than state, as the primitive units of computation. Formally, the system is defined over an event lattice $(E, \preceq)$, where each event represents an irreversible semantic action and the partial order encodes causal precedence. There is no operation that rewrites or deletes events. All system state is derived as a fold over $E$.

Let $\Sigma$ denote the semantic state space and let $f : \mathscr{P}(E) \to \Sigma$ be a deterministic reduction function. The current state at time $t$ is given by

$$\sigma_t = f(\{e \in E \mid e \preceq t\}).$$

Crucially, governance does not operate on $\sigma_t$ directly. It operates on admissibility relations over future events.

# 26 Formal Consent as an Admissibility Constraint

Within Spherepop OS, consent is not a preference but a constraint on the future evolution of the event lattice. A consent declaration introduces a forbidden region in semantic space that future events may not enter.

Let $\mathscr{C} \subset E$ be the set of consent events. Each $c \in \mathscr{C}$ induces a constraint function

$$\kappa_c : E \longrightarrow \{0, 1\},$$

where $\kappa_c(e) = 0$ indicates that event $e$ is inadmissible after $c$.

An event $e'$ is admissible at time $t$ if and only if

$$\forall c \in \mathscr{C}, \ c \preceq t \implies \kappa_c(e') = 1.$$

This construction has several immediate consequences. First, consent generalizes across instances. Blocking a semantic category blocks all future realizations, including generative variants. Second, consent is durable. It persists until explicitly revised by a new event. Third, consent is auditable. Violations are mechanically detectable as attempts to append inadmissible events.

Proposition 5: Any system in which consent is represented as mutable state rather than as an admissibility constraint is incapable of enforcing durable boundaries under optimization.

Proof: In a mutable-state system, consent may be overwritten by higher-priority updates or reinterpreted by downstream components without leaving a trace that invalidates future actions. Therefore, there exists a sequence of updates that preserves apparent compliance while violating user intent. This is structurally impossible in an event-constrained system, where admissibility is checked prior to state derivation.

# 27 Semantic Drift as a Consequence of Unconstrained State Evolution

Semantic drift arises when optimization operates over state representations whose meanings are not invariant under transformation. In recommender systems, content categories are treated as labels attached to items rather than as contracts governing behavior. As optimization pressures increase, producers exploit ambiguities in labeling to maximize engagement.

Spherepop OS prevents semantic drift by requiring that semantic transformations be explicit events. A category change is not an annotation but a transformation $\tau \in E$ that maps one semantic region to another. Such transformations are subject to the same admissibility constraints as all other events.

Let $\mathscr{S}$ denote the semantic manifold. Each event $e$ induces a transformation $T_e : \mathscr{S} \longrightarrow \mathscr{S}$. Semantic drift occurs when compositions of $T_e$ accumulate without constraint. Spherepop OS restricts admissible compositions by requiring that any transformation affecting category boundaries be explicitly declared and reviewable.

Proposition 6: Semantic drift is eliminated if and only if all category-altering transformations are explicit, irreversible events subject to consent constraints.

## 28 Why State-Based Optimization Collapses Under Acceleration

Acceleration exacerbates the failure of state-based systems because the gap between update speed and audit speed grows unbounded. When systems operate faster than humans can observe, mutable state becomes epistemically opaque. Even perfect logging is insufficient if interpretation lags action.

Event-sourced systems invert this relationship. Because admissibility is checked before events are committed, oversight operates ex ante rather than ex post. Verification scales with event rate, not state complexity.

This inversion is essential for surviving the transition to autonomous AI systems. Once AI agents can generate actions at superhuman speed, any governance model that relies on after-the-fact inspection is functionally void.

## 29 Polyxan and Typed Action Admissibility

Event-sourced semantics alone are insufficient to guarantee safety and coordination under accelerating intelligence. While they enforce ordering, auditability, and consent, they do not by themselves guarantee that proposed actions satisfy substantive correctness or safety properties. Polyxan addresses this limitation by introducing a typed action system in which admissibility is mediated by formal proof rather than trust or inference.

Polyxan models all actions as inhabitants of a type system whose types encode required invariants. Let $\mathscr{A}$ denote the space of action proposals and let $\mathscr{T}$ denote a set of action types parameterized by specifications. An action $a \in \mathscr{A}$ is admissible only if there exists a type $\tau \in \mathscr{T}$ and a proof object $\pi$ such that

$$\pi : a : \tau.$$

Here, the typing judgment simultaneously asserts that the action is well-formed and that it satisfies all invariants encoded in $\tau$. These invariants may include safety properties, consent constraints inherited from Spherepop OS, resource bounds, or coordination requirements.

Crucially, Polyxan separates the concerns of action generation and action verification. Agents may generate actions using any means, including opaque or stochastic processes. Governance operates solely on proof verification, which is deterministic, inspectable, and cheaper than generation.

## 30 Proof-Carrying Actions and Justified Confidence

The concept of justified confidence becomes precise in the Polyxan framework. Confidence in an action is justified if and only if its accompanying proof discharges all obligations imposed by the relevant specifications. Trust in agent intent or competence is neither required nor relevant.

Formally, let $\Phi$ be a predicate over actions representing a specification. An action $a$ is acceptable if there exists a proof $\pi$ such that

$$\pi \vdash \Phi(a).$$

Verification consists of checking $\pi$, not re-evaluating $a$ directly.

This structure yields an asymmetry essential for defense-favored systems. Proof construction is typically expensive and scales with action complexity. Proof verification is comparatively cheap and scales with proof size. As long as verification complexity grows sublinearly relative to generation complexity, oversight remains tractable under acceleration.

Proposition 7: In a proof-carrying system where verification cost is asymptotically lower than action generation cost, attackers do not gain advantage from increased generation speed.

Proof sketch: Any attempted exploit must be accompanied by a proof satisfying the same specifications as compliant actions. Faster generation does not reduce the cost of proof construction. Therefore increased speed yields no asymmetrical advantage.

## 31 Why Behavior-Based Alignment Is Strictly Weaker

Behavior-based alignment methods, including reinforcement learning from human feedback and preference inference, operate by adjusting agent policies to match observed human judgments. These methods do not provide guarantees; they provide statistical tendencies.

Let $\pi_\theta$ be a learned policy parameterized by $\theta$, trained to maximize expected reward $R$ derived from human feedback. Even if $\pi_\theta$ performs well on observed distributions, there exists no formal guarantee that it satisfies critical safety properties outside the training regime.

In contrast, Polyxan enforces properties universally through proof obligations. This distinction yields a strict dominance result.

Proposition 8: For any behavior-aligned system $S_b$, there exists a specification-enforced system $S_p$ such that $S_p$ satisfies all guarantees of $S_b$ and additional safety properties that $S_b$ cannot enforce.

Proof: Behavior-based alignment can be represented as optimization under a reward proxy. Any constraint enforced implicitly by the reward can be enforced explicitly as a specification in $S_p$. However, constraints requiring universal quantification or negative guarantees cannot be expressed reliably through reward signals alone. Therefore $S_b \subset S_p$ in expressive power.

This result does not imply that behavior-based methods are useless. It implies that they are incomplete as governance mechanisms. They may guide action generation, but they cannot substitute for verification.

## 32 Polyxan Versus Constitutional and Preference-Based AI

Recent alignment proposals emphasize constitutional AI, rule-following, or preference aggregation. These approaches attempt to internalize constraints within agent policies. Polyxan externalizes constraints instead.

Internalized constraints fail under distributional shift, adversarial pressure, or self-modification. Externalized constraints persist because they are enforced independently of agent internals. An agent may become more capable, faster, or differently motivated without affecting the validity of verification.

This distinction mirrors that between trusted code and verified code. Trusted code assumes compliance; verified code enforces it.

## 33    PlenumHub as a Coordination Substrate Under Adversarial Incentives

The final layer required to avoid chaos is coordination among multiple agents operating under competitive pressure. PlenumHub provides this layer by functioning as a shared verification and commitment environment.

PlenumHub is defined as a shared registry of events, proofs, and commitments. Agents interacting through PlenumHub must publish proposed actions along with proofs of admissibility. Other agents verify these proofs independently. Interaction is permitted only when all relevant policies are simultaneously satisfied.

Let $\mathscr{P}_i$ denote the policy specification of agent $i$. An interaction $a_{ij}$ between agents $i$ and $j$ is admissible if and only if there exists a proof $\pi$ such that

$$\pi \vdash \mathscr{P}_i(a_{ij}) \wedge \mathscr{P}_j(a_{ij}).$$

This construction eliminates unilateral advantage. Defection becomes visible and provable, enabling rapid response without trust assumptions.

## 34    Externality Encoding and Chaos Prevention

Chaos emerges when externalities are unpriced and untracked. PlenumHub addresses this by requiring that actions affecting shared resources include proofs of compliance with externality constraints.

Let $X$ denote a shared externality metric, such as risk exposure or environmental cost. Actions must satisfy bounds $X(a) \leq X_{max}$. These bounds are enforced mechanically rather than normatively.

Proposition 9: In a coordination system where all actions affecting shared resources are proof-constrained, tragedy-of-the-commons dynamics are structurally prevented.

Proof sketch: Tragedy-of-the-commons requires unilateral externalization without immediate consequence. Proof-constrained admissibility blocks such actions before they occur, regardless of actor intent.

## 35    Bidirectional Proof Topology and Navigable Governance

A critical but underexamined property of Polyxan is that it induces a bidirectional topology over actions and justifications. This property is not merely an implementation detail but a foundational

feature that transforms governance from a reactive process into a navigable semantic structure. In this respect, Polyxan departs sharply from conventional audit logs, compliance documents, and trust-based approval workflows, all of which are unidirectional and epistemically fragile.

In Polyxan, every admissible action is introduced into the system together with an explicit proof object that certifies its compliance with a specified set of invariants. This establishes a forward link from action to justification. Formally, if $a$ is an action and $\pi$ is a proof such that $\pi \vdash \Phi(a)$, then the system enforces the relation

$$a \rightarrow \pi.$$

This relation is total: no admissible action exists without an attached justification. The action is therefore never encountered as an opaque fact; it is always encountered as a claim accompanied by reasons.

The more consequential property, however, is the reverse relation. Proof objects in Polyxan are not ephemeral certificates but persistent semantic entities. Each proof induces a backward link to the set of actions whose admissibility depends upon it. Formally, for any proof $\pi$,

$$\pi \rightarrow \{a \in \mathcal{A} \mid \pi \vdash \Phi(a)\}.$$

This bidirectionality ensures that justifications are not merely explanatory but causally operative. If a proof is invalidated due to an error, a revoked assumption, or a superseded constraint, all dependent actions are automatically rendered inadmissible. Revocation propagates structurally rather than procedurally.

This topology contrasts fundamentally with prevailing governance mechanisms in AI systems. In most contemporary settings, actions are justified by reference to external documents, internal evaluations, or institutional authority. These references are informational rather than structural. An auditor may read a report asserting that a system was evaluated, but there is no mechanical pathway from the report back to the specific actions it was meant to authorize, nor any automatic invalidation if the report is later found deficient. Governance operates through narrative rather than linkage.

The Polyxan topology instead resembles systems in which meaning is constructed through reversible, inspectable operations. In interactive computational environments that emphasize transparency and learning, such as command-driven editors or text-based hypermedia systems, actions are intelligible precisely because they are embedded in a navigable structure of causes and effects. Commands can be inspected, reversed, recomposed, and understood in relation to their consequences. Polyxan generalizes this principle to governance itself. Actions are not merely performed; they are situated within a graph of justification that can be traversed, queried, and contested.

This property has direct implications for control under acceleration. As system capability increases, the number of actions proposed per unit time may grow beyond human comprehension. However, the cognitive burden on oversight does not scale with action volume but with proof complexity. Because proofs are explicit, modular, and reusable, the space of justification remains navigable even as action generation accelerates. Oversight is preserved not by slowing action but by preserving legibility through structure.

The bidirectional proof topology also resolves a central ambiguity in discussions of trust. Trust is replaced by inspectability. An agent or institution need not trust that another agent has acted responsibly; it can traverse the proof graph and verify compliance directly. Disagreement does not require speculation about intent; it is localized to specific proof obligations and assumptions.

This navigability property is essential for the subsequent analysis of multipolar coordination. Without bidirectional linkage, coordination collapses either into domination, where authority is centralized and opaque, or into chaos, where verification is too costly to sustain cooperation. Polyxan establishes the minimal epistemic substrate required for coordination without trust, which is a necessary precondition for avoiding both attractors simultaneously.

In the next section, this property is extended from single-system governance to multipolar environments via PlenumHub. The bidirectional topology that links actions to proofs within a system is generalized to link commitments, constraints, and externalities across systems, enabling coordination that remains stable even under adversarial incentives.

## 36 PlenumHub and the Geometry of Multipolar Power

The transition from single-system governance to multipolar coordination introduces a qualitatively different problem. When multiple agents or institutions possess significant autonomous capability, control failures no longer arise solely from internal opacity but from interaction effects between actors pursuing locally rational strategies. It is within this regime that the attractors of domination and chaos emerge.

PlenumHub is designed as a shared coordination substrate that generalizes the bidirectional proof topology of Polyxan across agents. Rather than functioning as a centralized authority, it operates as a common semantic and verification space in which commitments, constraints, and proofs are publicly instantiated and mechanically enforced.

Formally, let $\mathscr{A}_i$ denote the action space of agent $i$, and let $\mathscr{P}_i$ denote the policy specification governing that agent's admissible actions. PlenumHub introduces a shared interaction space $H$ such that an inter-agent action $a_{ij} \in \mathscr{A}_i \times \mathscr{A}_j$ is admissible if and only if there exists a proof $\pi$ satisfying

$$\pi \vdash \mathscr{P}_i(a_{ij}) \wedge \mathscr{P}_j(a_{ij}) \wedge \mathscr{P}_H(a_{ij}),$$

where $\mathscr{P}_H$ encodes shared invariants governing externalities, safety, and coordination norms.

This construction eliminates the need for bilateral trust. Agents need not assume benevolence or alignment in others; they need only verify proofs. Commitments become legible objects rather than informal promises, and violations become detectable events rather than contested narratives.

Importantly, PlenumHub does not require global consensus on values. It requires only agreement on enforceable invariants and verification procedures. This distinction is central to its ability to operate under adversarial incentives.

# 37 Domination and Chaos as Dynamical Attractors

Within multipolar systems, domination and chaos can be modeled as attractors in the phase space of power distribution and coordination capacity. Let $P(t)$ represent the concentration of effective power at time $t$, and let $C(t)$ represent the system's capacity for coordination, understood as the ability to enforce shared constraints.

Domination arises when $P(t)$ increases faster than $C(t)$ and becomes concentrated in a single actor or coalition. In this regime, oversight becomes impossible because the dominant actor can act beyond the temporal or epistemic reach of others. Chaos arises when $P(t)$ is distributed but $C(t)$ decays, forcing actors into defensive competition. Here, coordination collapses not due to malice but due to incentive incompatibility.

These dynamics can be expressed qualitatively as follows. Let $\dot{P}$ denote the rate of power accumulation and $\dot{C}$ the rate of coordination scaling. If

$$\dot{P} > \dot{C} \quad \text{and power is centralized,}$$

the system flows toward domination. If

$$\dot{P} > \dot{C} \quad \text{and power is distributed,}$$

the system flows toward chaos. In both cases, the inequality $\dot{P} > \dot{C}$ is decisive. Acceleration without constraint overwhelms governance.

Existing governance mechanisms attempt to counter these attractors through norms, treaties, or trust-building measures. However, such approaches scale poorly with capability growth because they rely on voluntary compliance and after-the-fact enforcement. As action speed increases, the lag between violation and response grows, rendering these mechanisms ineffective.

# 38 Formal Elimination of Domination and Chaos via Structural Constraints

PlenumHub, when combined with Spherepop OS and Polyxan, alters the dynamical landscape by changing the functional form of $C(t)$. Coordination capacity no longer depends primarily on human deliberation or institutional inertia, but on mechanical verification whose cost scales with proof complexity rather than action volume.

Let $V(t)$ denote the volume of actions proposed per unit time, and let $\kappa$ denote the average verification cost per action. In a proof-carrying system, total oversight cost scales as $\kappa V(t)$, while in a trust-based system, oversight cost scales superlinearly with $V(t)$ due to audit bottlenecks and dispute resolution.

Assume that proof verification is asymptotically cheaper than action generation, such that

$$\lim_{t \to \infty} \frac{\kappa V(t)}{P(t)} = 0.$$

Under this condition, coordination capacity scales with power rather than lagging behind it.

Proposition 10: In a multipolar system where all inter-agent actions are event-sourced, proof-carrying, and hub-coordinated, neither domination nor chaos is a stable attractor.

Proof: Domination requires unilateral action beyond contestation. However, unilateral action without proofs is inadmissible, and proof verification is independent of agent power. Therefore no actor can exceed coordination capacity by acceleration alone. Chaos requires incentive to defect by externalizing risk. However, externality-violating actions are blocked ex ante by shared invariants enforced through PlenumHub. Since defection yields no admissible payoff, competitive escalation does not dominate cooperative strategies. Thus both attractors are eliminated.

This result holds regardless of agent intent, relative capability, or internal architecture. It depends only on the existence of enforceable invariants and scalable verification. The narrow path is no longer a fragile balance but a structurally stable region enforced by system design.

## 39 Negative Externalities as First-Class Constraints

The defining feature of chaos-dominated regimes is not competition per se, but the systematic externalization of costs. Actors pursue locally rational strategies whose negative effects are displaced onto others, the environment, or the future. Traditional governance frameworks attempt to manage such externalities through regulation, liability, or moral suasion, all of which rely on delayed detection and discretionary enforcement. Under accelerating intelligence, these mechanisms fail because externalities propagate faster than they can be observed, attributed, or remedied.

The architectural approach developed in this paper treats negative externalities not as after-the-fact harms but as preconditions for admissibility. Externalities are elevated to first-class semantic objects that must be explicitly accounted for before action execution.

Formally, let $X = \{x_1, x_2, \ldots, x_n\}$ denote a vector of externality dimensions relevant to a shared system. These may include risk exposure, environmental impact, epistemic degradation, resource depletion, or coordination load. Each proposed action $a$ induces an externality vector

$$\mathbf{X}(a) = (x_1(a), x_2(a), \ldots, x_n(a)).$$

In conventional systems, $\mathbf{X}(a)$ is either unmeasured or estimated post hoc. In the Polyxan framework, admissibility requires a proof that $\mathbf{X}(a)$ lies within declared bounds. Let $\mathbf{B} = (b_1, b_2, \ldots, b_n)$ denote externally agreed limits. An action is admissible only if there exists a proof $\pi$ such that

$$\pi \vdash \forall i, \; x_i(a) \leq b_i.$$

These bounds are not global constants but context-sensitive invariants negotiated and encoded through PlenumHub. Different coordination domains may impose different externality budgets, reflecting heterogeneous risk tolerances without sacrificing enforceability.

# 40 Externality Budgets and Non-Local Accountability

A key failure mode of existing systems is that externalities are diffuse and non-local. An actor may reap immediate benefit while imposing small costs on many others, none of whom have standing to object individually. PlenumHub resolves this by making externality budgets collective and explicit.

Let $\mathscr{D}$ denote a coordination domain, such as a platform, ecosystem, or infrastructural layer. The domain maintains a cumulative externality ledger

$$L_{\mathscr{D}}(t) = \sum_{a \in \mathscr{A}_{\mathscr{D}}(t)} \mathbf{X}(a),$$

where $\mathscr{A}_{\mathscr{D}}(t)$ is the set of actions admitted up to time $t$.

Admissibility is conditional not only on per-action bounds but on aggregate constraints:

$$L_{\mathscr{D}}(t + 1) \leq \mathbf{B}_{\mathscr{D}}.$$

This construction prevents the incremental erosion characteristic of chaotic regimes. Even if each individual action appears harmless, cumulative harm is mechanically blocked once thresholds are approached. Importantly, enforcement does not require identifying a culprit or assigning blame; it requires only verifying ledger state.

# 41 Why This Eliminates the Incentive to Defect

Chaos emerges when defection yields higher expected payoff than cooperation. Defection is attractive precisely because externalities can be displaced. The proof-carrying externality model removes this advantage.

Proposition 11: In a system where all actions must satisfy explicit externality constraints ex ante, defection does not dominate cooperation.

Proof: Let defection be defined as an action that increases an actor's local payoff by externalizing cost. In a proof-constrained system, such an action is inadmissible unless the externality cost lies within allowed bounds. If bounds are exceeded, the action cannot occur. Therefore the payoff advantage of defection is eliminated at the admissibility layer, not compensated after the fact. Cooperation becomes the dominant strategy because it is the only strategy available.

This result is independent of actor intent, information asymmetry, or enforcement capacity. It follows directly from admissibility constraints.

# 42 Negative Externalities and the Domination Attractor

Domination arises when an actor can impose large externalities unilaterally, such as systemic risk or irreversible infrastructure changes. In traditional systems, such actions may be justified post hoc by authority or necessity. In the proposed stack, authority does not exempt actions from constraint.

Because externality proofs are required regardless of actor identity, power does not grant exemption. A dominant actor attempting to deploy a high-risk action must still produce a proof that the action lies within agreed bounds. If such a proof does not exist, the action is structurally blocked.

Thus domination is eliminated not by distributing power evenly, but by equalizing constraint. Capability may differ; admissibility does not.

## 43 Externalities, Consent, and Semantic Pollution

Many of the harms discussed earlier, such as unwanted content exposure and semantic drift, can be understood as externalities imposed on users' cognitive environments. Treating these as first-class externalities allows consent to be enforced mechanically.

In Spherepop OS, semantic pollution is represented as an increase in entropy within a declared semantic region. Consent events define acceptable entropy bounds for a user's informational environment. Actions that would exceed these bounds, such as repeated exposure to excluded categories or manipulative content patterns, are inadmissible regardless of engagement metrics.

This reframes content governance as an externality management problem rather than a moderation problem. The goal is not to judge content but to prevent unconsented cognitive load.

## 44 Externality Accounting as the Core of the Narrow Path

The narrow path between domination and chaos is ultimately a path of bounded externalities. Domination corresponds to unbounded externality imposition by a single actor. Chaos corresponds to unbounded cumulative externalities imposed by many actors. Both are eliminated when externalities are constrained at the point of action.

The Spherepop–Polyxan–PlenumHub stack enforces this constraint uniformly. Externalities are declared, measured, bounded, and verified before action occurs. This shifts governance from reaction to prevention and from moral judgment to structural design.

In the next section, this externality-centric view is used to compare the proposed architecture with existing multipolar governance proposals, highlighting why most fail to address externalities at the correct level of abstraction.

## 45 Comparison with Existing Multipolar Governance Proposals

Existing approaches to multipolar AI governance largely fall into three categories: norm-based coordination, treaty-based regulation, and trust- or reputation-based institutional oversight. While these approaches differ in emphasis, they share a common structural weakness. All operate downstream of action, relying on voluntary compliance, delayed detection, or discretionary enforcement. Under conditions of accelerating intelligence, this temporal lag becomes fatal.

Norm-based proposals assume that shared values and expectations can stabilize behavior through mutual restraint. Such frameworks function tolerably well when violations are rare, visible, and slow-

moving. However, when actions are generated at machine speed and effects propagate nonlocally, norms lose their coordinating power. Actors who violate norms gain immediate advantage, while punishment, if it arrives at all, arrives too late to prevent harm. Norms therefore degrade into symbolic commitments rather than operative constraints.

Treaty-based governance attempts to formalize coordination through legally binding agreements. While treaties can encode explicit limits, they depend on human interpretation, enforcement, and adjudication. These processes scale poorly with action volume and are vulnerable to ambiguity, loopholes, and strategic reinterpretation. Moreover, treaties typically bind states rather than systems, leaving open the question of how autonomous agents and private actors are constrained in practice.

Reputation-based and institutional oversight models emphasize transparency, audits, and accountability mechanisms. These approaches improve visibility but do not change admissibility conditions. An action may be audited, criticized, or penalized after execution, yet the system provides no structural means to prevent recurrence. As capability increases, oversight becomes increasingly ceremonial, producing reports rather than control.

By contrast, the architecture developed in this paper operates at a different layer of abstraction. Instead of attempting to align incentives ex post, it enforces constraints ex ante. Instead of relying on trust, it relies on verification. Instead of regulating actors, it regulates actions. This shift explains why existing proposals struggle to address domination and chaos simultaneously. They attempt to manage outcomes without constraining the processes that generate them.

Importantly, the proposed stack does not require global agreement on values, political systems, or moral doctrines. It requires only agreement on formal invariants and verification procedures. This minimalism is what allows coordination to scale across heterogeneous actors without collapsing into either centralized authority or anarchic competition.

# 46 Unifying the Stack as a Plenum and Field Model

The full significance of the Spherepop–Polyxan–PlenumHub stack becomes apparent when viewed as a unified plenum rather than as a collection of discrete components. A plenum, in the classical sense, is a continuous medium that mediates interaction, propagation, and constraint. In this context, the plenum is not physical but semantic and operational: a structured field in which actions, proofs, and externalities propagate according to invariant laws.

Spherepop OS defines the temporal and semantic geometry of this plenum. Events are points in the field, ordered by causality and irreversibility. Meaning arises not from instantaneous state but from trajectories through event space. Consent defines forbidden regions within this field, carving out volumes that actions may not enter.

Polyxan defines the local admissibility conditions within the plenum. Proof objects function as field tensors that encode whether a proposed action satisfies the invariants at its point of insertion. An action without a valid proof is analogous to a particle violating conservation laws; it cannot exist within the field.

PlenumHub defines the global coupling between regions of the field. It ensures that interactions

across domains respect shared constraints and that externalities propagate through explicit channels rather than leaking invisibly. Coordination emerges not from negotiation at every interaction but from shared field structure.

Formally, let $\mathscr{F}$ denote the plenum, and let actions be represented as localized excitations $a \in \mathscr{F}$. Each excitation carries an externality vector $\mathbf{X}(a)$ and a proof field $\pi(a)$. The dynamics of the system are governed by admissibility equations of the form

$$\mathscr{L}(\mathscr{F}, a) = 0,$$

where $\mathscr{L}$ encodes semantic, safety, consent, and coordination invariants. Actions that violate these equations are not evolved; they are excluded.

This field-theoretic perspective clarifies why domination and chaos are eliminated structurally. Domination corresponds to singularities in which one excitation overwhelms the field, bypassing invariants. Chaos corresponds to turbulence in which unbounded excitations accumulate externalities faster than constraints can dissipate them. In the proposed plenum, neither singularities nor unbounded turbulence are admissible solutions. The field equations forbid them.

The narrow path is therefore not a balancing act performed by agents but a property of the field itself. As long as invariants are enforced, trajectories remain within the stable region. Capability growth corresponds to higher-energy excitations, not to violation of conservation laws.

## 47 Alignment Reconsidered

This unification allows a final reframing of alignment. Alignment is not primarily a property of agent preferences, reward functions, or internal representations. It is a property of the environment in which agents act. An aligned system is one in which misaligned actions are structurally inadmissible.

From this perspective, the alignment problem dissolves into an engineering problem: how to design a plenum in which agency, consent, and coordination remain conserved quantities even as intelligence accelerates. The Spherepop–Polyxan–PlenumHub stack offers one such design. It does not promise utopia, consensus, or moral convergence. It promises something more modest and more fundamental: that no actor, human or artificial, can unilaterally destroy the conditions under which choice, cooperation, and meaning remain possible.

## 48 Closing Synthesis

The narrow path between domination and chaos is real, but it is not navigated by vigilance or virtue. It is navigated by structure. Domination arises when power escapes constraint. Chaos arises when constraint fails to scale. Both failures are consequences of architectures that treat consent, accountability, and externalities as secondary concerns.

By elevating these concepts to first-class invariants enforced through event sourcing, proof-carrying admissibility, and shared coordination substrates, it is possible to design systems in which accelera-

tion does not imply loss of control. The future of artificial intelligence need not be decided by trust or fear. It can be decided by design.

# A  Appendix A: Formal Grammar of Event-Sourced and Proof-Carrying Systems

This appendix presents a precise syntactic grammar for event-sourced, proof-carrying coordination systems of the kind analyzed in the main text. The grammar is intentionally conservative. It specifies the minimal structural requirements for well-formedness while remaining agnostic about semantic interpretation, proof theory, or execution strategy. Its purpose is to define a common syntactic substrate upon which the categorical and sheaf-theoretic semantics developed in subsequent appendices may be imposed.

The grammar is expressed in Backus–Naur Form (BNF). Terminal symbols are written in typewriter font, non-terminal symbols are enclosed in angle brackets, and metavariables range over syntactic classes unless otherwise specified.

## A.1  Lexical Primitives and Atomic Types

We assume a standard lexical layer sufficient to identify symbolic entities, scalar values, and Boolean constants. These primitives are intentionally minimal, as higher-level structure is enforced compositionally rather than lexically.

```
<Identifier> ::= <Letter> { <Letter> | <Digit> | "_" }

<Letter>     ::= "A" | ... | "Z" | "a" | ... | "z"
<Digit>      ::= "0" | ... | "9"

<Scalar>     ::= <Integer> | <Float>
<Integer>    ::= <Digit> { <Digit> }
<Float>      ::= <Digit> { <Digit> } "." <Digit> { <Digit> }

<Bool>       ::= "true" | "false"
```

Identifiers are treated as globally unique names within a given namespace. No syntactic distinction is made between identifiers denoting agents, actions, predicates, or specifications; such distinctions are imposed semantically in later appendices.

## A.2  Composite Value Structures

The grammar supports structured values for representing collections, vectors, and symbolic references. These structures allow externality vectors, argument lists, and constraint payloads to be expressed uniformly.

```
<Value> ::= <Scalar>
          | <Bool>
```

```
              | <Identifier>
              | <Set>
              | <Vector>

<Set>    ::= "{" [ <ValueList> ] "}"
<Vector> ::= "[" [ <ValueList> ] "]"

<ValueList> ::= <Value>
              | <Value> "," <ValueList>
```

No ordering semantics are imposed on sets, while vectors are assumed to be ordered. The grammar does not enforce type homogeneity within collections; typing constraints are deferred to the specification layer.

## A.3   Event Declarations

All system evolution is expressed through events. Events are the sole source of state change and are irreducible. There is no notion of mutable state outside the event stream.

```
<Event> ::= "event" <Identifier> ":" <EventType> <EventPayload>
```

Each event has a unique identifier, a declared event type, and an associated payload.

## A.4   Event Types

The grammar distinguishes event classes only to enforce minimal well-formedness conditions. The semantics of these classes are defined structurally rather than procedurally.

```
<EventType> ::= "consent"
              | "action"
              | "transformation"
              | "proof"
              | "revocation"
```

The inclusion of a distinct `revocation` event type is essential. Revocation is not treated as a negation of prior events but as a first-class event whose effects propagate through dependency structure.

## A.5   Event Payload Structure

Event payloads are structured records containing declarations. Declarations bind identifiers to values within the local scope of the event.

```
<EventPayload> ::= "{" <DeclarationList> "}"

<DeclarationList> ::= <Declaration>
                    | <Declaration> "," <DeclarationList>

<Declaration> ::= <Identifier> "=" <Value>
```

No fixed schema is imposed on payloads at the grammar level. Schema enforcement is handled by specification predicates and proofs.

## A.6 Consent and Constraint Specifications

Consent is represented syntactically as a declarative constraint on semantic regions. These regions are hierarchical and compositional.

```
<ConsentSpec> ::= "forbid" <SemanticRegion>
                | "allow" <SemanticRegion>

<SemanticRegion> ::= <Identifier>
                   | <Identifier> "::" <SemanticRegion>
```

The use of a recursive region grammar enforces prefix-closure of semantic scopes. A constraint on a region applies transitively to all subregions unless overridden by a more specific declaration.

## A.7 Action Declarations

Actions are parameterized symbolic constructs. They do not imply execution; they are proposals subject to admissibility judgment.

```
<Action> ::= "action" <Identifier> "(" [ <ArgumentList> ] ")"

<ArgumentList> ::= <Argument>
                 | <Argument> "," <ArgumentList>

<Argument> ::= <Identifier> ":" <Type>
```

The grammar intentionally separates action identity from admissibility. An action declaration alone carries no authority.

## A.8 Type and Specification Syntax

Types are treated syntactically as identifiers. Specifications are expressed as conjunctions of predicates applied to arguments.

```
<Type> ::= <Identifier>

<Proof> ::= "proof" <Identifier> ":" <Specification>

<Specification> ::= <Predicate>
                  | <Predicate> "AND" <Specification>

<Predicate> ::= <Identifier> "(" [ <ArgumentList> ] ")"
```

Specifications are purely propositional at the syntactic level. The grammar does not privilege any particular logic; it merely requires explicit enumeration of obligations.

### A.9  Admissibility Judgments

Admissibility is the sole syntactic bridge between actions and proofs. No action is considered well-formed for execution unless it appears within an admissibility judgment.

```
<Admissible> ::= <Action> "⊢" <Proof>
```

An admissibility judgment is a syntactic claim that a given proof discharges all obligations required by the corresponding action. The grammar does not verify this claim; it only enforces that the claim is explicit and structurally present.

### A.10  Syntactic Invariants

The grammar enforces the following invariants:

First, all system evolution is event-based; there is no syntax for implicit state mutation. Second, all authority is mediated through proofs; there is no syntax for unverified action. Third, revocation is explicit and compositional; there is no silent erasure. Fourth, consent is declarative and hierarchical; there is no ad hoc filtering.

These invariants are necessary but not sufficient for correctness. Their sufficiency is established only when the grammar is interpreted through the categorical and sheaf-theoretic semantics developed in Appendices B through D.

# B  Appendix B: Categorical Formulation of Events, Proofs, and Admissibility

This appendix provides a categorical formalization of the event-sourced, proof-carrying coordination system defined syntactically in Appendix A. The goal is not abstraction for its own sake, but the precise identification of the compositional structures that ensure scalability, auditability, and stability

under acceleration. Category theory is employed because it allows invariants to be expressed independently of implementation details, while making explicit the conditions under which composition is well-defined.

Throughout, we work in the setting of small categories. All constructions are internal to this setting unless otherwise noted.

## B.1   The Event Category

Let $\mathscr{E}$ be a category whose objects are events as defined syntactically in Appendix A. Morphisms in $\mathscr{E}$ represent causal precedence.

For events $e_1, e_2 \in \mathrm{Ob}(\mathscr{E})$, there exists a morphism

$$e_1 \rightarrow e_2$$

if and only if event $e_1$ causally precedes event $e_2$.

We impose the following axioms:

First, identity: every event precedes itself. Second, transitivity: causal precedence composes. Third, antisymmetry: if $e_1 \rightarrow e_2$ and $e_2 \rightarrow e_1$, then $e_1 = e_2$.

These axioms imply that $\mathscr{E}$ is a thin category, and therefore equivalent to a partially ordered set. The absence of nontrivial automorphisms enforces irreversibility: events cannot be undone, only followed by further events.

Acyclicity is enforced by the absence of non-identity endomorphisms, ensuring that no event can be causally downstream of itself except trivially. This categorical property formalizes the event-sourced nature of the system.

## B.2   The Action Category

Let $\mathscr{A}$ be a category whose objects are action declarations. Morphisms in $\mathscr{A}$ represent lawful transformations of actions, such as refinement, parameter instantiation, or sequential composition, provided such transformations preserve required specifications.

Composition in $\mathscr{A}$ is partial: two actions may be composable only if the output context of the first is compatible with the input context of the second. This partiality reflects the fact that not all actions may be sequenced meaningfully.

The identity morphism on an action corresponds to the trivial refinement of that action.

## B.3   The Proof Category

Let $\mathscr{P}$ be a category whose objects are proof artifacts. A proof artifact is a structured justification discharging one or more specifications. Morphisms in $\mathscr{P}$ represent proof refinement, extension, or weakening, provided that discharged specifications remain satisfied.

Unlike $\mathscr{A}$, the category $\mathscr{P}$ is generally not thin. Multiple distinct proofs may discharge the same specification, and morphisms between proofs capture relationships such as proof reuse, abstraction,

or specialization.

## B.4 The Specification Category

Let $\mathcal{S}$ be a category whose objects are specifications. A specification is understood as a formal obligation or invariant that may be required of an action and discharged by a proof.

Morphisms in $\mathcal{S}$ represent implication or refinement relations between specifications. If there exists a morphism

$$\phi_1 \rightarrow \phi_2,$$

then satisfaction of $\phi_1$ implies satisfaction of $\phi_2$.

This ordering allows specifications to be composed and compared structurally.

## B.5 Functorial Structure

There exist two canonical functors connecting these categories.

First, a functor

$$H \,:\, \mathcal{A} \rightarrow \mathcal{S}$$

that assigns to each action the specification it requires. This functor formalizes the notion that actions impose obligations on the system.

Second, a functor

$$G \,:\, \mathcal{P} \rightarrow \mathcal{S}$$

that assigns to each proof the specification it discharges. This functor formalizes the notion that proofs certify satisfaction of obligations.

These functors are required to be compatible with composition: refining an action refines its required specification, and refining a proof refines the specification it discharges.

## B.6 Admissibility as a Pullback

Admissible actions are defined categorically as the pullback of $H$ and $G$ over $\mathcal{S}$. Concretely, we define a category $\mathcal{A}_{\mathrm{adm}}$ by the pullback square

$$\mathcal{A}_{\mathrm{adm}} = \mathcal{A} \times_{\mathcal{S}} \mathcal{P}.$$

An object of $\mathcal{A}_{\mathrm{adm}}$ is a pair $(a, p)$ such that

$$H(a) = G(p).$$

That is, the specification required by the action is exactly the specification discharged by the proof.

Morphisms in $\mathcal{A}_{\mathrm{adm}}$ are pairs of morphisms in $\mathcal{A}$ and $\mathcal{P}$ that commute with the pullback condition. This ensures that refinement of admissible actions preserves admissibility.

## B.7 Compositionality of Admissible Actions

The pullback construction guarantees compositional closure. If $(a_1, p_1)$ and $(a_2, p_2)$ are admissible and composable, then their composition $(a_2 \circ a_1, p_2 \circ p_1)$ is admissible provided the composed proof discharges the composed specification.

This property is critical under acceleration. As action volume increases, admissibility checking remains local and compositional rather than global and procedural.

## B.8 Faithful Embedding of Behavior-Based Systems

Let $\mathscr{B} \subset \mathscr{A}$ be the full subcategory of behavior-aligned actions, that is, actions whose correctness is inferred post hoc from observed behavior rather than proven ex ante.

There exists a faithful functor

$$I : \mathscr{B} \hookrightarrow \mathscr{A}_{\text{adm}},$$

which embeds behavior-based actions into the proof-carrying system by pairing them with trivial or opaque proofs.

However, there exists no functor

$$\mathscr{A}_{\text{adm}} \rightarrow \mathscr{B}$$

that preserves safety invariants. This asymmetry establishes a strict expressive dominance: proof-carrying governance can represent all behavior-based systems, but not vice versa.

## B.9 Implications for Control and Scaling

The categorical formulation makes explicit why proof-carrying systems scale under acceleration. Admissibility is not a global property computed over entire histories; it is a local property enforced by pullback structure. As a result, increases in action throughput do not increase the complexity of admissibility checking beyond the complexity of individual proofs.

This structural property underlies the main text's claim that domination and chaos are eliminated not by monitoring outcomes, but by constraining composition itself.

# C   Appendix C: Sheaf-Theoretic Interpretation of Consent, Externalities, and Coordination

This appendix develops a sheaf-theoretic semantics for the event-sourced, proof-carrying system introduced in Appendices A and B. The purpose of this interpretation is to make precise how local admissibility conditions compose into global coherence, and how failures of such composition correspond exactly to negative externalities, coordination breakdowns, and pathological attractors.

Sheaf theory is employed because it provides a mathematically precise language for reasoning about locality, compatibility, and gluing. These notions are central to the system under consideration:

actions are proposed locally, consent is expressed contextually, and coordination requires that locally admissible actions remain globally consistent.

## C.1   The Semantic Base Space

Let $X$ be a topological space whose points represent semantic contexts. A semantic context is defined as a minimal situation in which the admissibility of an action can be evaluated. Examples include a specific user's content environment, a particular subsystem of an infrastructure, or a localized interaction domain within a multipolar coordination setting.

The topology on $X$ encodes refinement of context. An open set $U \subseteq X$ represents a region of semantic coherence: a set of contexts within which the same consent constraints, specifications, and admissibility judgments apply. If $V \subseteq U$, then $V$ is a refinement of $U$, corresponding to a more specific or constrained context.

The topological structure is not assumed to be Hausdorff or metrizable. Only the axioms necessary for sheaf theory are required. In particular, the existence of arbitrary unions and finite intersections of open sets is sufficient.

## C.2   Presheaf of Admissible Actions

Define a presheaf

$$\mathscr{F} \; : \; \mathcal{O}(X)^{\mathrm{op}} \longrightarrow \mathbf{Set}$$

on the category of open sets of $X$. For each open set $U \subseteq X$, the set $\mathscr{F}(U)$ consists of all actions that are admissible within semantic region $U$.

Restriction maps

$$\rho_{U,V} \; : \; \mathscr{F}(U) \longrightarrow \mathscr{F}(V) \quad \text{for } V \subseteq U$$

model contextual refinement. An action admissible in a broader context may or may not remain admissible when additional constraints are imposed.

The presheaf axioms enforce that restriction is identity-preserving and compositional.

## C.3   Consent as Subsheaf Construction

Consent declarations impose exclusions on admissible actions. Formally, consent is represented by a subsheaf

$$\mathscr{F}_{\mathrm{consent}} \subseteq \mathscr{F}.$$

For each open set $U$, the set $\mathscr{F}_{\mathrm{consent}}(U)$ consists of those actions in $\mathscr{F}(U)$ that do not violate any consent constraint applicable within $U$.

The subsheaf condition ensures that consent constraints are monotone under refinement: if an action is forbidden in a broader context, it remains forbidden in all refined contexts. Conversely, allowances may be introduced only by moving to larger contexts or by explicit revocation events.

This formalizes the durability and non-locality of consent boundaries emphasized in the main text.

## C.4   Externalities as Failures of Gluing

The defining property of a sheaf is the gluing condition. Given an open cover $\{U_i\}$ of an open set $U$, a family of sections

$$a_i \in \mathscr{F}(U_i)$$

that agree on all overlaps $U_i \cap U_j$ must glue uniquely to a section $a \in \mathscr{F}(U)$.

In the present framework, negative externalities correspond exactly to failures of this gluing condition. A family of locally admissible actions may fail to glue because their combined effects violate a global invariant, such as an externality budget, safety constraint, or coordination requirement.

Such failures are not semantic accidents; they are structural obstructions. The obstruction is global even if no local violation is visible.

## C.5   Cocycle Conditions and Cumulative Harm

More formally, consider a cover $\{U_i\}$ and a family $\{a_i\}$ of locally admissible actions. Compatibility on overlaps expresses pairwise non-interference. However, cumulative harm arises when higher-order compatibility fails: triple overlaps, or the total effect over the union $U$, violate a constraint not expressible at the level of any single $U_i$.

This failure corresponds to a nontrivial cocycle in the Čech cohomology of the presheaf. The cocycle represents a latent externality that is invisible locally but unavoidable globally.

Thus negative externalities are identified with cohomological obstructions to global sections.

## C.6   PlenumHub as a Global Section Enforcer

PlenumHub is interpreted sheaf-theoretically as enforcing the sheaf condition by construction. Only those families of local actions that admit a global section are executable.

Formally, PlenumHub restricts the system to the sheafification

$$\mathscr{F}^+$$

of the presheaf $\mathscr{F}$. Actions that remain presheaf-local but fail to sheafify are rejected prior to execution.

This guarantees that coordination is not emergent or negotiated post hoc. It is enforced by the topology of admissibility itself.

## C.7   Domination and Chaos in Sheaf-Theoretic Terms

Domination corresponds to singular global sections that ignore local restrictions. In sheaf-theoretic language, these are sections that exist only because restriction maps are bypassed or disabled. Such

sections cannot exist in a true sheaf.

Chaos corresponds to families of locally admissible sections that cannot be glued due to unresolved cocycle obstructions. In the absence of sheaf enforcement, such families accumulate and interact destructively.

By enforcing sheaf conditions mechanically, the system excludes both pathologies. Only globally coherent sections are admitted, and all admissible local actions must participate in a consistent global structure.

## C.8   Stability Under Refinement

A crucial consequence of the sheaf-theoretic formulation is stability under refinement. As semantic contexts become more detailed or constraints become more specific, the admissible action space shrinks monotonically, but coherence is preserved.

This property underlies the system's ability to remain stable under increasing intelligence and complexity. Refinement does not introduce inconsistency; it merely excludes previously admissible trajectories that are no longer coherent under stricter constraints.

# D   Appendix D: Structural Elimination of Pathological Attractors

This appendix completes the formal foundation by demonstrating, at the level of system structure rather than agent behavior, the elimination of the two pathological attractors analyzed in the main text: domination and chaos. The argument synthesizes the syntactic constraints of Appendix A, the categorical compositionality of Appendix B, and the sheaf-theoretic coherence conditions of Appendix C.

The central claim is that domination and chaos are not merely unlikely under the proposed architecture; they are structurally inadmissible. They correspond to objects or trajectories that cannot exist within the formal system defined by admissibility, composition, and gluing.

## D.1   Attractors as Structural Fixed Points

We begin by formalizing the notion of an attractor in the present context. Let $\mathscr{T}$ denote the space of admissible trajectories, where a trajectory is a finite or infinite sequence of events drawn from the event category $\mathscr{E}$ and whose actions are objects of the admissibility category $\mathscr{A}_{\mathrm{adm}}$.

An attractor is defined as a subset $A \subseteq \mathscr{T}$ such that, under refinement of context, increase in capability, or extension of history, trajectories generically evolve toward $A$ unless prevented by explicit constraints.

The claim of this appendix is that, in the proposed system, the subsets corresponding to domination and chaos are empty. There exist no admissible trajectories exhibiting either property.

## D.2 Domination as Structural Singularity

Domination is characterized by the unilateral execution of actions whose effects cannot be contested, audited, or constrained by other actors or by the system itself. Structurally, domination requires the existence of at least one of the following:

First, an action executed without an associated proof discharging its required specifications. Second, an action whose proof cannot be inspected or verified by other participants. Third, an action that bypasses consent or externality constraints due to privileged status of the actor.

Each of these conditions corresponds to a violation of structural invariants.

By Appendix A, no action is well-formed for execution unless it appears in an admissibility judgment pairing it with a proof. By Appendix B, admissibility is defined categorically as a pullback; there is no morphism admitting an action without a corresponding proof object. By Appendix C, global execution requires compatibility with all relevant local constraints; singular sections are excluded.

Therefore, domination would require the existence of a morphism in $\mathscr{A}_{\mathrm{adm}}$ whose projection to $\mathscr{P}$ is undefined or opaque. No such morphism exists.

## D.3 Lemma: Non-Existence of Uncertified Actions

Formally, let $a \in \mathrm{Ob}(\mathscr{A})$ be an action. Suppose there exists a trajectory in which $a$ is executed. Then there exists a proof object $p \in \mathrm{Ob}(\mathscr{P})$ such that $(a, p) \in \mathrm{Ob}(\mathscr{A}_{\mathrm{adm}})$.

Proof. This follows directly from the definition of $\mathscr{A}_{\mathrm{adm}}$ as a pullback. Objects of $\mathscr{A}_{\mathrm{adm}}$ are precisely such pairs. There is no inclusion of $\mathscr{A}$ into the execution semantics independent of $\mathscr{A}_{\mathrm{adm}}$. □

Thus domination, understood as action beyond contestation, is structurally impossible.

## D.4 Chaos as Failure of Global Coherence

Chaos is characterized by a multipolar regime in which actors follow locally admissible strategies whose aggregate effects undermine shared invariants. Structurally, chaos requires the existence of a family of locally admissible actions whose combined execution violates a global constraint, such as an externality bound or safety invariant.

In the sheaf-theoretic formulation of Appendix C, this corresponds precisely to a failure of the gluing condition. Locally admissible sections exist, but they do not glue to a global section.

## D.5 Lemma: Non-Existence of Non-Glueable Executable Families

Let $\{U_i\}$ be an open cover of a semantic region $U \subseteq X$, and let $a_i \in \mathscr{F}(U_i)$ be locally admissible actions. Suppose the family $\{a_i\}$ fails to glue to a global section in $\mathscr{F}(U)$. Then the execution of $\{a_i\}$ is inadmissible.

Proof. By construction, PlenumHub restricts execution to the sheafification $\mathscr{F}^+$. Only families admitting a global section are executable. Non-glueable families are rejected ex ante. □

Thus chaos, understood as cumulative incoherence, is excluded structurally rather than mitigated dynamically.

## D.6  Externalities and the Closure of the Argument

The essential bridge between domination and chaos is the treatment of negative externalities. Domination corresponds to unbounded externality imposition by a single actor; chaos corresponds to bounded local externalities whose accumulation is unbounded globally.

By requiring explicit externality accounting at the level of admissibility, and by enforcing gluing conditions across contexts, both modes of failure are eliminated simultaneously.

No actor may impose unbounded externalities unilaterally, because admissibility requires proof against global invariants. No collection of actors may impose unbounded externalities cumulatively, because non-glueable action families are rejected.

## D.7  Stability Under Acceleration

A critical concern addressed in the main text is whether these guarantees persist as system capability accelerates. Structurally, acceleration increases the density of proposed actions but does not alter the admissibility criteria.

Because admissibility is local, compositional, and proof-based, the cost of verification scales with proof complexity rather than with action volume. The categorical and sheaf-theoretic constructions remain invariant under increased throughput.

Thus neither domination nor chaos can re-emerge as attractors under acceleration. They are not unstable equilibria awaiting perturbation; they are absent states in the system's phase space.

## D.8  Conclusion

The elimination of domination and chaos is achieved not through balance, negotiation, or trust, but through structure. By enforcing event-sourced irreversibility, proof-carrying admissibility, categorical compositionality, and sheaf-theoretic coherence, the system defines a space of admissible trajectories from which pathological attractors are excluded by construction.

This completes the formal foundation underlying the arguments of the main text. The final appendices give formalizations in lean.

# E  Appendix E: Lean Formalizations of Event-Sourced Admissibility, Proof-Carrying Actions, and Externality Constraints

This appendix provides Lean formalizations corresponding to the core structures used in the manuscript: event-sourced histories, consent as admissibility constraints, proof-carrying actions, and externality-bounded coordination. The code is written to be readable and modular, using standard Lean 4 idioms. It is intended as a schematic formalization suitable for refinement into a fully mechanized development.

## E.1 E.1 Preliminaries and Basic Types

```
namespace Plenum

universe u v

-- Abstract identifiers for events, actions, agents, and specifications.
constant Event   : Type
constant Action  : Type
constant Agent   : Type
constant Spec    : Type
constant ExtDim  : Type    -- Externality dimension (e.g., risk, pollution, epist

-- A proof object that a specification holds of an action.
-- We treat Spec as a predicate on Action.
abbrev Holds (▯ : Action → Prop) (a : Action) : Prop := ▯ a

-- Externality vector: for each action and externality dimension, a scalar cost.
constant Cost     : Type
constant leCost   : Cost → Cost → Prop
infix:50 " ▯ₖ " => leCost

-- Axioms needed to reason about costs.
axiom le_refl  : ∀ c : Cost, c ▯ₖ c
axiom le_trans : ∀ {a b c : Cost}, a ▯ₖ b → b ▯ₖ c → a ▯ₖ c
```

## E.2 E.2 Event Order, Histories, and Event-Sourced State

We model event-sourcing by a causal precedence relation on events and histories as finite lists of
events. The key property is that admissibility depends on history, not merely on a mutable state flag.

```
-- Causal precedence (a partial order) on events.
constant precedes : Event → Event → Prop
infix:50 " ▯ " => precedes

axiom precedes_refl   : ∀ e : Event, e ▯ e
axiom precedes_trans  : ∀ {e1 e2 e3 : Event}, e1 ▯ e2 → e2 ▯ e3 → e1 ▯ e3
axiom precedes_antisymm : ∀ {e1 e2 : Event}, e1 ▯ e2 → e2 ▯ e1 → e1 = e2

-- Histories as finite event lists.
abbrev History := List Event
```

```
-- Membership predicate for histories.
def InHist (e : Event) (h : History) : Prop := e ∈ h
```

### E.3    E.3 Consent as an Admissibility Constraint

Consent is formalized as a predicate on actions parameterized by history. A consent event induces a
constraint that restricts admissible future actions.

```
-- A consent constraint: given a history, determine if an action is allowed.
abbrev Consent := History → Action → Prop


-- A family of consent constraints; in practice, there may be many.
abbrev ConsentSet := List Consent


-- All consent constraints must approve an action for it to be admissible.
def ConsentOK (Cs : ConsentSet) (h : History) (a : Action) : Prop :=
  ∀ C ∈ Cs, C h a


-- Monotonicity: once forbidden by consent, remains forbidden under history exte
-- This captures durability of consent boundaries under event-sourcing.
def MonotoneConsent (C : Consent) : Prop :=
  ∀ (h₁ h₂ : History) (a : Action),
    (h₁ <:+ h₂) → (C h₁ a → C h₂ a)


-- A common strengthening: if an action is disallowed at some point, extension ca
def PersistentBan (C : Consent) : Prop :=
  ∀ (h₁ h₂ : History) (a : Action),
    (h₁ <:+ h₂) → (¬ C h₁ a → ¬ C h₂ a)
```

### E.4    E.4 Proof-Carrying Actions (Polyxan Core Judgment)

Polyxan admissibility requires that an action satisfy both consent constraints and a set of formal spec-
ifications. We represent specifications as predicates on actions and proofs as Lean terms inhabiting
those predicates.

```
-- A "specification" is simply a predicate on actions.
abbrev Specification := Action → Prop
abbrev SpecSet := List Specification


-- An action is SpecOK if it satisfies every specification in the set.
def SpecOK (⊠s : SpecSet) (a : Action) : Prop :=
  ∀ ⊠ ∈ ⊠s, ⊠ a
```

```
-- Proof-carrying admissibility: consent + specs.
def Admissible (Cs : ConsentSet) (⬚s : SpecSet) (h : History) (a : Action) : Prop
   ConsentOK Cs h a ⬚ SpecOK ⬚s a

-- A "certificate" packages an action with proofs of admissibility.
structure Certificate (Cs : ConsentSet) (⬚s : SpecSet) (h : History) where
   act   : Action
   ok    : Admissible Cs ⬚s h act
```

## E.5   E.5 Bidirectional Proof Topology as Dependency Tracking

We capture the bidirectional link property by representing a proof object as a reusable lemma and
defining a dependency relation from proofs to certified actions.

```
-- A reusable proof artifact is modeled as a theorem/lemma object.
-- For formal dependency tracking, we model it as a predicate that implies SpecOK

structure ProofArtifact (⬚s : SpecSet) where
  witness : ∀ a : Action, SpecOK ⬚s a → True  -- Placeholder for a richer artifact

-- Forward link: any certificate contains a proof (Lean term) of Admissible.
def forwardLink {Cs ⬚s h} (cert : Certificate Cs ⬚s h) : Admissible Cs ⬚s h cert.a
   cert.ok

-- Backward link: given a set of certificates, we can collect those that depend o
-- In Lean, dependency is syntactic; here we model it semantically.
def DependsOn (⬚s : SpecSet) (a : Action) : Prop := SpecOK ⬚s a

def BackwardImage (⬚s : SpecSet) (A : Set Action) : Set Action :=
   {a | a ∈ A ⬚ DependsOn ⬚s a}
```

This semantic abstraction is adequate for the manuscript's claim: revocation corresponds to changing $\Phi$ or its trusted axioms so that previously derivable proofs no longer typecheck, thereby invalidating dependent certificates.

## E.6   E.6 Externality Vectors and Budgets

Externalities are modeled as cost functions per dimension with per-action and cumulative constraints.

```
-- Externality measure: for each action and dimension, a cost.
constant extCost : Action → ExtDim → Cost
```

```
-- A per-dimension budget
abbrev Budget := ExtDim → Cost


-- An action respects the budget if all its dimensional costs are bounded.
def ExtOK (B : Budget) (a : Action) : Prop :=
  ∀ d : ExtDim, extCost a d ⊑ₖ B d


-- Add an externality constraint to the specification set.
def ExtSpec (B : Budget) : Specification := fun a => ExtOK B a
```

For cumulative constraints, we model a ledger and require that each admitted action preserves global boundedness.

```
-- A ledger is a mapping from dimension to cumulative cost.
abbrev Ledger := ExtDim → Cost


-- Ledger update rule after admitting an action.
constant addCost : Cost → Cost → Cost
axiom add_monotone_left : ∀ {x y z : Cost}, x ⊑ₖ y → addCost x z ⊑ₖ addCost y z
axiom add_monotone_right : ∀ {x y z : Cost}, x ⊑ₖ y → addCost z x ⊑ₖ addCost z y


def ledgerUpdate (L : Ledger) (a : Action) : Ledger :=
  fun d => addCost (L d) (extCost a d)


-- Global budget condition on a ledger.
def LedgerOK (B : Budget) (L : Ledger) : Prop :=
  ∀ d : ExtDim, L d ⊑ₖ B d


-- An action is globally admissible w.r.t. ledger if it preserves LedgerOK.
def PreservesLedger (B : Budget) (L : Ledger) (a : Action) : Prop :=
  LedgerOK B L → LedgerOK B (ledgerUpdate L a)
```

### E.7  E.7 PlenumHub-Style Multipolar Interaction as Joint Policy Satisfaction

We model a hub-mediated interaction between two agents as an action that must satisfy both agents' policies and hub policies. Policies are specifications parameterized by agent.

```
-- Agent policy: a specification predicate on actions.
abbrev Policy := Agent → Specification


-- Hub policy: a specification predicate (not agent-indexed).
```

```
abbrev HubPolicy := Specification

-- Joint admissibility for interaction action between agents i and j.
def JointAdmissible
    (Pi : Policy) (Pj : Policy) (Ph : HubPolicy)
    (a : Action) (i j : Agent) : Prop :=
    (Pi i) a □ (Pj j) a □ Ph a
```

## E.8    E.8 Structural Prevention of Defection via Externality Constraints

We formalize the key claim: actions that attempt to externalize beyond bounds are inadmissible, there-
fore the defection advantage is eliminated at the admissibility layer.

```
-- A defection action is one that violates externality budget.
def Defection (B : Budget) (a : Action) : Prop := ¬ ExtOK B a

theorem defection_inadmissible
  (Cs : ConsentSet) (□s : SpecSet) (h : History) (B : Budget) (a : Action) :
  (ExtSpec B ∈ □s) → Defection B a → ¬ Admissible Cs □s h a := by
  intro hExt hDef
  unfold Admissible
  intro hAdm
  have hSpecOK : SpecOK □s a := hAdm.right
  have hExtOK : ExtSpec B a := by
    have : (ExtSpec B) ∈ □s := hExt
    exact hSpecOK (ExtSpec B) this
  have : ExtOK B a := hExtOK
  exact hDef this
```

## E.9    E.9 Domination and Chaos as Non-Admissible Trajectory Properties

A full mechanization of domination and chaos requires a system dynamics model.  Here we state
the core structural form: trajectories are sequences of certified actions; inadmissible actions cannot
appear.  Domination corresponds to unilateral constraint bypass; chaos corresponds to cumulative
externality overflow. Both are excluded by construction.

```
-- A trajectory is a sequence of certified actions under evolving histories.
structure Step (Cs : ConsentSet) (□s : SpecSet) where
  h    : History
  cert : Certificate Cs □s h

abbrev Trajectory (Cs : ConsentSet) (□s : SpecSet) := List (Step Cs □s)
```

```
-- A domination attempt would require an action that violates hub policy or bypas
-- Such an action cannot be packaged as a Certificate.
-- Chaos corresponds to cumulative ledger overflow; prevented by requiring Prese

-- Placeholder predicates capturing these forbidden properties.
constant DominationAttempt : Action → Prop
constant ChaosOverflow      : Action → Prop

theorem domination_excluded_by_certification
  (Cs : ConsentSet) (⊠s : SpecSet) (h : History) (a : Action) :
 DominationAttempt a → ¬ (∃ cert : Certificate Cs ⊠s h, cert.act = a) := by
  intro hDom
 -- In a complete development, DominationAttempt implies failure of some spec i
 -- Therefore no Admissible proof exists, hence no Certificate exists.
 -- This theorem is a schematic statement, to be refined with a concrete spec en
  intro hex
 cases hex with
 | intro cert hEq =>
  -- Replace with contradiction once DominationAttempt is tied to specificatio
   cases hDom
```

### E.10    E.10 Notes on Mechanization Strategy

A complete mechanization proceeds by (i) instantiating specifications as concrete predicates, (ii) en-coding consent constraints as monotone admissibility functions over histories, (iii) defining hub poli-cies to include externality budgets and safety invariants, and (iv) proving that trajectories built from certificates cannot realize domination or chaos properties. The key is that the relevant invariants are stated as specifications and therefore become proof obligations, rather than being inferred from behavior.

    end Plenum

## F    Appendix F: Lean Development with Operational Semantics, Ordered Monoids, Trajectory Validity, and Inductive Safety Proofs

This appendix refines Appendix E by replacing schematic placeholders with a concrete, inductive development. We formalize (i) an ordered commutative monoid of costs, (ii) a ledger update semantics, (iii) trajectory validity as an inductive predicate, and (iv) theorems showing that externality overflow and constraint-bypass actions cannot occur in valid trajectories. The development remains abstract in

the sense that it does not commit to a particular SMT backend or proof-carrying code format; rather, it makes precise the structural logic that the manuscript relies upon.

## F.1 F.1 Ordered Commutative Monoid of Costs

We require costs to form a commutative monoid equipped with a preorder compatible with addition. This is the minimal algebraic structure for externality accounting.

```
namespace Plenum2

universe u

constant Cost  : Type
constant le    : Cost → Cost → Prop
infix:50 " ≤ₖ " => le


-- Commutative monoid structure
constant zero  : Cost
constant add   : Cost → Cost → Cost
infixl:65 " ⊕ " => add


-- Preorder laws
axiom le_refl  : ∀ c : Cost, c ≤ₖ c
axiom le_trans : ∀ {a b c : Cost}, a ≤ₖ b → b ≤ₖ c → a ≤ₖ c


-- Monoid laws
axiom add_assoc : ∀ a b c : Cost, (a ⊕ b) ⊕ c = a ⊕ (b ⊕ c)
axiom add_comm  : ∀ a b : Cost, a ⊕ b = b ⊕ a
axiom add_zero  : ∀ a : Cost, a ⊕ zero = a


-- Order compatibility (monotonicity)
axiom add_mono_left  : ∀ {a b c : Cost}, a ≤ₖ b → (a ⊕ c) ≤ₖ (b ⊕ c)
axiom add_mono_right : ∀ {a b c : Cost}, a ≤ₖ b → (c ⊕ a) ≤ₖ (c ⊕ b)
```

## F.2 F.2 Actions, Specifications, Consent, and Proof-Carrying Certificates

We model specifications as predicates on actions, consent as a history-indexed predicate, and admissibility as the conjunction of consent and specification satisfaction.

```
constant Event  : Type
constant Action : Type
```

```
abbrev History := List Event

abbrev Specification := Action → Prop
abbrev SpecSet := List Specification

abbrev Consent := History → Action → Prop
abbrev ConsentSet := List Consent

def SpecOK (⬚s : SpecSet) (a : Action) : Prop :=
  ∀ ⬚ ∈ ⬚s, ⬚ a

def ConsentOK (Cs : ConsentSet) (h : History) (a : Action) : Prop :=
  ∀ C ∈ Cs, C h a

def Admissible (Cs : ConsentSet) (⬚s : SpecSet) (h : History) (a : Action) : Prop
  ConsentOK Cs h a ⬚ SpecOK ⬚s a

structure Certificate (Cs : ConsentSet) (⬚s : SpecSet) (h : History) where
  act : Action
  ok  : Admissible Cs ⬚s h act
```

## F.3   F.3 Externality Ledger Semantics and Budget Constraints

We define externality dimensions, a cost function per dimension, a ledger as a function from dimensions to costs, and an update rule that accumulates costs additively.

```
constant ExtDim : Type
constant extCost : Action → ExtDim → Cost

abbrev Budget := ExtDim → Cost
abbrev Ledger := ExtDim → Cost

def LedgerOK (B : Budget) (L : Ledger) : Prop :=
  ∀ d : ExtDim, L d ⬚_k B d

def ledgerUpdate (L : Ledger) (a : Action) : Ledger :=
  fun d => (L d) ⬚ (extCost a d)

def ExtOK (B : Budget) (a : Action) : Prop :=
  ∀ d : ExtDim, extCost a d ⬚_k B d
```

```
-- Two externality styles: per-action boundedness and cumulative ledger boundedr
-- The manuscript's stronger claim is obtained via cumulative boundedness.

def PreservesLedger (B : Budget) (L : Ledger) (a : Action) : Prop :=
  LedgerOK B L → LedgerOK B (ledgerUpdate L a)
```

### F.4  F.4 Policy Encoding: Hub Policy Requires Ledger Preservation

PlenumHub is modeled as requiring that each admitted action preserve a shared externality ledger bound. This is a structural externality internalization constraint.

```
def HubPolicy (B : Budget) (L : Ledger) : Specification :=
  fun a => PreservesLedger B L a
```

### F.5  F.5 Operational Semantics of Steps and Trajectory Validity

A trajectory is a sequence of steps. Each step carries a history, a ledger state, and a certificate whose admissibility must include the hub's ledger-preservation policy at that step. Validity is defined inductively, forcing the ledger to remain within bounds at every step.

```
structure Step (Cs : ConsentSet) (⬛s : SpecSet) where
  h    : History
  L    : Ledger
  cert : Certificate Cs ⬛s h

abbrev Trajectory (Cs : ConsentSet) (⬛s : SpecSet) := List (Step Cs ⬛s)

-- A predicate stating that a step's certificate includes the hub policy at ledge
 def StepHasHubPolicy (Cs : ConsentSet) (⬛s : SpecSet) (B : Budget) (s : Step Cs ⬛s
  (HubPolicy B s.L) ∈ ⬛s

-- Inductive validity: start from an initial ledger within budget; each step must
-- (i) have a ledger within budget, (ii) include the hub policy, and
-- (iii) preserve the ledger to the next step.

inductive ValidTrajectory (Cs : ConsentSet) (⬛s : SpecSet) (B : Budget) : Ledger
| nil (L0 : Ledger) :
    LedgerOK B L0 →
    ValidTrajectory Cs ⬛s B L0 []
| cons (L0 : Ledger) (s : Step Cs ⬛s) (ts : Trajectory Cs ⬛s) :
    LedgerOK B L0 →
    s.L = L0 →
```

```
      StepHasHubPolicy Cs ⊠s B s →
 -- the certificate must entail hub policy satisfaction for the action
   (let a := s.cert.act
    let okSpecs := s.cert.ok.right
    (HubPolicy B L0) a) →
   -- next ledger is the updated ledger
   ValidTrajectory Cs ⊠s B (ledgerUpdate L0 s.cert.act) ts →
   ValidTrajectory Cs ⊠s B L0 (s :: ts)
```

The above definition makes the hub constraint causally operative: ledger preservation is checked at each step, and the next state is defined by the ledger update rule. Valid trajectories therefore cannot accumulate externalities beyond the declared budget.

## F.6   F.6 Main Theorem I: No Externality Overflow in Valid Trajectories

We formalize overflow as the existence of a step whose ledger violates the budget. The theorem shows that overflow cannot occur.

```
def Overflow (B : Budget) (L : Ledger) : Prop := ¬ LedgerOK B L

theorem no_overflow_in_valid
  (Cs : ConsentSet) (⊠s : SpecSet) (B : Budget) :
  ∀ (L0 : Ledger) (ts : Trajectory Cs ⊠s),
    ValidTrajectory Cs ⊠s B L0 ts →
    ∀ s ∈ ts, ¬ Overflow B s.L := by
  intro L0 ts hValid
  induction hValid with
  | nil L0 hOK =>
      intro s hsIn
      cases hsIn
  | cons L0 s ts hOK hEqL stepHas hubSat hTail ih =>
      intro s' hsIn
      cases hsIn with
      | head =>
          -- s' = s
          subst_vars
          unfold Overflow
          intro hOv
          exact hOv (by
            -- s.L = L0 and LedgerOK B L0
            simpa [hEqL] using hOK)
```

```
        | tail hsInTail =>
            -- apply IH to tail
            have : ¬ Overflow B s'.L := ih s' hsInTail
            exact this
```

This theorem is the formal core of the manuscript's externality claim: cumulative harm is prevented because ledger boundedness is an invariant of valid trajectories.

## F.7   F.7 Constraint Bypass and the Impossibility of Uncertified Actions

Domination, in its structural form, requires bypassing the admissibility gate. In the formal system, bypass corresponds to an action appearing in the evolution without a certificate. We model a trace as the list of actions extracted from a trajectory and prove that every action in a valid trajectory has an admissibility proof.

```
def ActionsOf (Cs : ConsentSet) (⬚s : SpecSet) (ts : Trajectory Cs ⬚s) : List Act
  ts.map (fun s => s.cert.act)

theorem every_action_certified
  (Cs : ConsentSet) (⬚s : SpecSet) (B : Budget) :
  ∀ (L0 : Ledger) (ts : Trajectory Cs ⬚s),
    ValidTrajectory Cs ⬚s B L0 ts →
    ∀ a ∈ ActionsOf Cs ⬚s ts,
      ∃ (h : History) (L : Ledger), ∃ (cert : Certificate Cs ⬚s h),
        cert.act = a ⬚ Admissible Cs ⬚s h a := by
  intro L0 ts hValid
  induction hValid with
  | nil L0 hOK =>
      intro a haIn
      cases haIn
  | cons L0 s ts hOK hEqL stepHas hubSat hTail ih =>
      intro a haIn
      simp [ActionsOf] at haIn
      cases haIn with
      | inl hHead =>
          -- a is the head action
          refine ⬚s.h, s.L, s.cert, ?_, ?_⬚
          · simpa [hHead]
          · -- admissible follows from certificate
            -- the certificate stores Admissible Cs ⬚s s.h a
          have hAdm : Admissible Cs ⬚s s.h s.cert.act := s.cert.ok
```

```
        simpa [hHead] using hAdm
   | inr hTailIn =>
        exact ih a hTailIn
```

This theorem expresses the non-bypass property: any realized action is necessarily backed by a proof of admissibility. Structural domination requires an actor to act without being subject to constraints; this is ruled out by the construction of trajectories.

## F.8 F.8 Eliminating the Defection Advantage with Ledger Preservation

To connect back to the chaos attractor, we show that actions that would cause immediate ledger overflow cannot be admitted at any valid step. This is stronger than a per-action bound; it is a cumulative guarantee.

```
def CausesOverflowNext (B : Budget) (L : Ledger) (a : Action) : Prop :=
  Overflow B (ledgerUpdate L a)


theorem no_overflow_causing_action_admitted
  (Cs : ConsentSet) (⬚s : SpecSet) (B : Budget) :
  ∀ (L0 : Ledger) (s : Step Cs ⬚s),
    LedgerOK B L0 →
    s.L = L0 →
    StepHasHubPolicy Cs ⬚s B s →
    (let a := s.cert.act; (HubPolicy B L0) a) →
    ¬ CausesOverflowNext B L0 s.cert.act := by
  intro L0 s hOK hEqL hHas hubSat
  unfold CausesOverflowNext Overflow
  intro hOv
 -- By hub policy satisfaction, ledgerUpdate preserves LedgerOK, contradiction
  have hPres : PreservesLedger B L0 s.cert.act := hubSat
  have : LedgerOK B (ledgerUpdate L0 s.cert.act) := hPres hOK
  exact hOv this
```

This theorem formalizes the manuscript's key externality mechanism: defection by externality overflow is not merely punishable; it is inadmissible.

## F.9 F.9 Notes on Completing the Domination/Chaos Attractor Proof in Lean

The manuscript's attractor elimination result can be fully mechanized by defining domination and chaos as properties of trajectories. A domination attempt is modeled as the existence of an action in the evolution that is not certified or that violates a non-negotiable hub invariant. A chaos regime is modeled as the existence of a trajectory whose ledger exceeds budget or whose policies become

inconsistent. The theorems above already provide the two core lemmas needed: (i) all actions are certified, hence bypass is impossible, and (ii) ledger overflow cannot occur in valid trajectories, hence cumulative externality escalation is structurally prevented.

To capture multipolar interaction explicitly, one extends `Action` to include agent parameters and composes agent policies conjunctively within `SpecSet`. The proofs proceed by induction on `ValidTrajectory` and rely only on monotonicity and compositionality of specifications. The critical point remains unchanged: the system's evolution is restricted to trajectories whose steps are admissible, and admissibility is defined in terms of proof-checkable invariants.

end Plenum2

# G   Appendix G: Runnable Lean 4 Formalization with Concrete Costs

This appendix provides a concrete, executable Lean 4 development instantiating the abstract framework of the previous appendices. Costs are natural numbers, externalities accumulate additively, and all invariants are proved using Lean's standard library. This appendix demonstrates that the elimination of domination and chaos is not merely schematic but mechanically verified.

```
import Std

namespace PlenumConcrete

open Std


----------------------------------------------------------------
-- G.1 Basic Types
----------------------------------------------------------------

abbrev Event   := Nat
abbrev Action := Nat
abbrev History := List Event


----------------------------------------------------------------
-- G.2 Costs and Externalities
----------------------------------------------------------------

abbrev Cost  := Nat
abbrev ExtDim := Nat

abbrev Budget := ExtDim → Cost
abbrev Ledger := ExtDim → Cost
```

```
def extCost (a : Action) (d : ExtDim) : Cost :=
  (a + d) % 3    -- simple bounded cost model


----------------------------------------------------------------
-- G.3 Ledger Semantics
----------------------------------------------------------------

def LedgerOK (B : Budget) (L : Ledger) : Prop :=
  ∀ d, L d ⧸ B d

def ledgerUpdate (L : Ledger) (a : Action) : Ledger :=
  fun d => L d + extCost a d

def PreservesLedger (B : Budget) (L : Ledger) (a : Action) : Prop :=
  LedgerOK B L → LedgerOK B (ledgerUpdate L a)


----------------------------------------------------------------
-- G.4 Specifications, Consent, and Admissibility
----------------------------------------------------------------

abbrev Specification := Action → Prop
abbrev SpecSet := List Specification

abbrev Consent := History → Action → Prop
abbrev ConsentSet := List Consent

def SpecOK (⧸s : SpecSet) (a : Action) : Prop :=
  ∀ ⧸ ∈ ⧸s, ⧸ a

def ConsentOK (Cs : ConsentSet) (h : History) (a : Action) : Prop :=
  ∀ C ∈ Cs, C h a

def Admissible (Cs : ConsentSet) (⧸s : SpecSet) (h : History) (a : Action) : Prop
  ConsentOK Cs h a ⧸ SpecOK ⧸s a

structure Certificate (Cs : ConsentSet) (⧸s : SpecSet) (h : History) where
  act : Action
  ok  : Admissible Cs ⧸s h act
```

56

```
-----------------------------------------------------------------
-- G.5 Hub Policy (Externality Constraint)
-----------------------------------------------------------------


def HubPolicy (B : Budget) (L : Ledger) : Specification :=
  fun a => PreservesLedger B L a


-----------------------------------------------------------------
-- G.6 Operational Semantics
-----------------------------------------------------------------


structure Step (Cs : ConsentSet) (⬚s : SpecSet) where
  h    : History
  L    : Ledger
  cert : Certificate Cs ⬚s h

abbrev Trajectory (Cs : ConsentSet) (⬚s : SpecSet) := List (Step Cs ⬚s)


-----------------------------------------------------------------
-- G.7 Valid Trajectories (Inductive)
-----------------------------------------------------------------


inductive ValidTrajectory
  (Cs : ConsentSet) (⬚s : SpecSet) (B : Budget) :
  Ledger → Trajectory Cs ⬚s → Prop
| nil (L0 : Ledger) :
    LedgerOK B L0 →
    ValidTrajectory L0 []
| cons (L0 : Ledger) (s : Step Cs ⬚s) (ts : Trajectory Cs ⬚s) :
    LedgerOK B L0 →
    s.L = L0 →
    HubPolicy B L0 ∈ ⬚s →
    (HubPolicy B L0) s.cert.act →
    ValidTrajectory (ledgerUpdate L0 s.cert.act) ts →
    ValidTrajectory L0 (s :: ts)


-----------------------------------------------------------------
-- G.8 No Externality Overflow Theorem
-----------------------------------------------------------------
```

```
def Overflow (B : Budget) (L : Ledger) : Prop :=
  ¬ LedgerOK B L

theorem no_overflow
  (Cs : ConsentSet) (⬚s : SpecSet) (B : Budget) :
  ∀ (L0 : Ledger) (ts : Trajectory Cs ⬚s),
    ValidTrajectory Cs ⬚s B L0 ts →
    ∀ s ∈ ts, ¬ Overflow B s.L := by
  intro L0 ts hV
  induction hV with
  | nil L0 hOK =>
      intro s hIn
      cases hIn
  | cons L0 s ts hOK hEq hHub hSat hTail ih =>
      intro s' hIn
      cases hIn with
      | head =>
          subst_vars
          unfold Overflow
          intro hOv
          exact hOv hOK
      | tail hTailIn =>
          exact ih s' hTailIn


-------------------------------------------------------------
-- G.9 All Actions Are Certified (No Domination)
-------------------------------------------------------------

def ActionsOf (Cs : ConsentSet) (⬚s : SpecSet) (ts : Trajectory Cs ⬚s) : List Act
  ts.map (fun s => s.cert.act)

theorem every_action_certified
  (Cs : ConsentSet) (⬚s : SpecSet) (B : Budget) :
  ∀ (L0 : Ledger) (ts : Trajectory Cs ⬚s),
    ValidTrajectory Cs ⬚s B L0 ts →
    ∀ a ∈ ActionsOf Cs ⬚s ts,
      ∃ h L cert,
        cert.act = a ⬚ Admissible Cs ⬚s h a := by
  intro L0 ts hV
  induction hV with
```

```
    | nil L0 hOK =>
        intro a hIn
        cases hIn
    | cons L0 s ts hOK hEq hHub hSat hTail ih =>
        intro a hIn
        simp [ActionsOf] at hIn
        cases hIn with
        | inl hHead =>
            refine ⟨s.h, s.L, s.cert, ?_, ?_⟩
            · simpa [hHead]
            · simpa [hHead] using s.cert.ok
        | inr hTailIn =>
            exact ih a hTailIn


---------------------------------------------------------------
-- G.10 Interpretation
---------------------------------------------------------------


-- Any trajectory that evolves according to ValidTrajectory
-- (1) cannot overflow the externality budget (no chaos),
-- (2) cannot contain uncertified actions (no domination).
-- Both properties are invariants proved by induction.

 end PlenumConcrete
```

# References

[1] Ammann, N. (2024). *How to Avoid Two AI Catastrophes: Domination and Chaos*. Talk and interview series on AI governance, high-assurance systems, and multipolar coordination. Transcripts and recordings available via public alignment research archives.

[2] National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research. (1979). *The Belmont Report: Ethical Principles and Guidelines for the Protection of Human Subjects of Research*. U.S. Department of Health, Education, and Welfare.

[3] Floridi, L. (2013). *The Ethics of Information*. Oxford University Press.

[4] Floridi, L. (2016). *Onlife Manifesto: Being Human in a Hyperconnected Era*. Springer.

[5] Owen, R., Macnaghten, P., & Stilgoe, J. (2013). Responsible research and innovation: From science in society to science for society, with society. *Science and Public Policy*, 39(6), 751–760.

[6] Yudkowsky, E. (2008). Artificial intelligence as a positive and negative factor in global risk. In N. Bostrom & M. Ćirković (Eds.), *Global Catastrophic Risks*. Oxford University Press.

[7] Christiano, P., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2018). Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*.

[8] Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*.

[9] Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.

[10] Kant, I. (1785). *Groundwork of the Metaphysics of Morals*. Translated editions, Cambridge University Press.

[11] Mill, J. S. (1859). *On Liberty*. John W. Parker and Son.

[12] Zuboff, S. (2019). *The Age of Surveillance Capitalism*. PublicAffairs.

[13] Pariser, E. (2011). *The Filter Bubble: What the Internet Is Hiding from You*. Penguin Press.

[14] Appel, A. W. (1997). Foundational proof-carrying code. In *Proceedings of the 12th Annual IEEE Symposium on Logic in Computer Science*.

[15] Necula, G. C. (1997). Proof-carrying code. In *Proceedings of the 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*.

[16] Lamport, L. (1998). The part-time parliament. *ACM Transactions on Computer Systems*, 16(2), 133–169.

[17] Ostrom, E. (1990). *Governing the Commons: The Evolution of Institutions for Collective Action*. Cambridge University Press.

[18] Varian, H. R. (2019). Artificial intelligence, economics, and industrial organization. In *The Economics of Artificial Intelligence*. University of Chicago Press.

[19] Tufekci, Z. (2015). Algorithmic harms beyond Facebook and Google. *Colorado Technology Law Journal*, 13, 203–218.

[20] Kroll, J. A., Huey, J., Barocas, S., et al. (2017). Accountable algorithms. *University of Pennsylvania Law Review*, 165, 633–705.

[21] European Union. (2016). *General Data Protection Regulation* (EU) 2016/679.

[22] European Union. (2022). *Digital Services Act*.

[23] Simon, H. A. (1969). *The Sciences of the Artificial*. MIT Press.

[24] Goodhart, C. A. E. (1975). Problems of monetary management: The UK experience. In *Papers in Monetary Economics*. Reserve Bank of Australia.