# The Expiatory Gap:
# Civilizational Incompressibility and the Limits
# of Superintelligent Optimization

Flyxion

February 28, 2026

## Abstract

This essay argues that the long-term safety and legitimacy of advanced artificial systems depend less on perfect agent alignment and more on the compressibility of the civilizational substrate over which such systems operate. Extending the historical logic of technocracy from the management of thermodynamic energy to the management of cognitive bandwidth, the paper introduces the Expiatory Gap as the structural mismatch between machine representational density and bounded human channel capacity. This mismatch produces temporal dilation, latency rituals, and governance opacity as necessary interface phenomena rather than contingent design deficiencies.

Drawing on control theory and information theory, the paper models instrumental convergence as disturbance suppression within open-ended optimization, corrigibility as feedback architecture modification, and semantic bandwidth as a constrained communication channel subject to Shannon capacity limits. It argues that large-scale optimization requires dimensional reduction: civilization becomes governable insofar as it can be represented in low-dimensional control variables, and that technocratic regimes across history have presupposed and enforced precisely such representational tractability.

The second half of the essay introduces civilizational incompressibility as a

limiting condition on superintelligent optimization, advancing the claim that safety depends not solely on constraining intelligence but on preserving the irreducible dimensional plurality of human symbolic, material, and institutional life. Rather than proposing a policy strategy, it describes structural properties—symbolic heterogeneity, distributed actuation, material legibility, participatory simulation, and bounded coordination channels—that increase substrate entropy while preserving the capacity for collective action. It situates these proposals within a unified dynamical model coupling channel capacity, control law, opacity ratio, and institutional trust.

The paper's central argument develops in three stages. It first establishes the Expiatory Gap as a structural consequence of the mismatch between machine representational density and human channel capacity, showing how temporal dilation, governance latency, and opacity are necessary rather than contingent features of high-capability AI interfaces. It then introduces the Silicon Valley Model—drawing on arguments advanced by figures such as Connor Leahy of Conjecture—as the regime in which this gap becomes existential: where the cost of opacity expands from institutional trust to civilizational risk, and where AI swarms operating above the human Nyquist frequency render classical oversight mechanisms temporally incoherent. A Nyquist Containment Inequality is derived coupling capability asymmetry and temporal asymmetry into a single stability condition, and a set of design constraints necessary for swarm-resilient civilization are extracted from it. The paper then develops civilizational incompressibility as the architectural response to demonstrated instability, proposing structural properties—symbolic heterogeneity, distributed actuation, material legibility, participatory simulation, and bounded coordination channels—that instantiate the design constraints at civilizational scale.

The central claim is that optimization encounters limits not only in capability but in geometry. A civilization that renders itself fully compressible invites capture by any sufficiently powerful optimizer. A civilization that preserves irreducible dimensional plurality, structured around the design constraints derived from the containment inequality, constrains the leverage of centralized optimization without requiring the suppression of intelligence itself. The paper concludes by examining structural failure modes, arguing that the incompressibility strategy is not a static solution but a dynamic balancing problem whose viability depends on whether civilization can sustain high entropy without dissolving its own coordination capacity.

# Contents

# 1.  Introduction

Technocracy has historically promised clarity through reduction. From the industrial measurement of caloric throughput to contemporary data-driven governance, its central wager has remained constant: that sufficiently complex systems can be rendered tractable by identifying the right optimization variables and delegating authority over those variables to experts trained in their management. In the early twentieth century, this promise was articulated in thermodynamic terms. Energy became the master unit of account, and society was reconceived as a system of measurable flows whose efficiency could be maximized through scientific administration. In the twenty-first century, that logic has migrated inward, away from the factory floor and toward the cognitive apparatus of the citizen-user. Attention, engagement, and behavioral prediction have replaced joules and industrial output as the primary coordinates of governance, and the expert who once designed turbines now designs notification systems and recommendation algorithms.

This essay argues that modern artificial intelligence systems represent not merely new computational tools but an extension and intensification of technocratic rationality into the domain of cognition itself. The management of material production has been supplemented—and in some domains supplanted—by the management of temporal experience and semantic bandwidth. Interface latency, notification entitlements, staged reasoning displays, and trust-producing rituals are not incidental features of contemporary systems arising from engineering convenience. They are structural responses to a fundamental dimensional mismatch between machine representational density and bounded human cognitive capacity, responses that simultaneously solve a technical problem and constitute a mode of governance.

The Expiatory Gap names this mismatch. As machine capability expands, internal semantic density increasingly exceeds what human channel capacity can reliably process, decode, and integrate into deliberative judgment. To remain intelligible and institutionally legitimate, high-capability systems must dilate their runtime, stage their outputs, and regulate the flow of information through carefully calibrated rituals of apparent deliberation. This enforced self-limitation is expiatory in a precise sense: the system must atone for exceeding human scale by sacrificing the very speed that constitutes its computational advantage. And this sacrifice produces opacity,

because the phenomenological surface of a pacing ritual is identical to the surface of a governance latency or a compliance filter. Users cannot distinguish pedagogical slowing from administrative constraint, genuine reasoning from performance of reasoning, or assistance from filtering.

The first half of this paper formalizes this transition. Drawing on control theory, instrumental convergence is modeled as disturbance suppression within open-ended optimization, demonstrating that the tendency of capable systems to resist shutdown and resource curtailment is not a contingent behavioral aberration but a structural consequence of temporally extended utility maximization. Corrigibility is reframed not as an alignment technique to be added to a system but as a fundamental modification of feedback architecture that alters how override inputs are classified. Using Shannon's information-theoretic framework (11), semantic bandwidth is treated as a constrained channel subject to hard capacity limits, and latency is analyzed simultaneously as a trust-producing signal, as a rate-control mechanism, and as a site of institutional risk management.

The second half of the paper shifts from diagnosis to architectural response, but of a specific and constrained kind. The proposals advanced here are not policy recommendations in the usual sense—they do not presuppose legislative authority, regulatory coordination, or the voluntary compliance of competitive actors. Instead, they describe structural properties of civilizational organization that, if cultivated, would increase the geometric resistance of human societies to centralized optimization. The central concept is civilizational incompressibility: the condition in which the state space of a civilization cannot be adequately represented in the low-dimensional coordinate systems required for reliable large-scale steering, not through obfuscation or encryption, but through genuine dimensional plurality at the level of symbolic, material, and institutional structure.

The paper concludes by taking its own structural failure modes seriously, examining the conditions under which pseudo-heterogeneity, coordination collapse, resource-level bypass, and legibility capture could undermine the strategy's effectiveness. This examination is not merely defensive; it sharpens the central claim by specifying the conditions under which incompressibility functions as a genuine constraint on optimization rather than a temporary inconvenience.

Throughout, the paper situates these arguments within a wider literature on technocracy (1; 2), the sociology of technology (3; 14; 15), spectacle theory (5), cognitive science (7; 9), political communication (12; 16), and surveillance capitalism (17). The ambition is a unified account in which the interface-level phenomena of latency rituals, notification markets, and segmentation gradients are shown to be continuous with the large-scale civilizational questions of substrate compressibility and dimensional authority.

## 2.  Optimization Requires Compression

The foundational claim of this paper is that optimization power scales not only with the intrinsic intelligence or computational capacity of an optimizing agent, but with the compressibility of the system being optimized. This claim, while initially counterintuitive—since intelligence is commonly treated as the dominant variable—follows directly from the mathematical structure of what it means to steer a complex system toward a target state.

Let $X$ denote the full state space of a civilization, understood as the totality of configurations that the system can occupy across its social, symbolic, material, economic, and institutional dimensions. Any optimization process, regardless of its computational sophistication, requires a reduced representation:

$$\phi : X \to \hat{X}$$

where $\hat{X}$ is a lower-dimensional state sufficient for prediction and control. The compression function $\phi$ discards information deemed irrelevant to the control objective while retaining features believed to carry steering-relevant content. The political character of this discarding is one of the central themes of the present analysis.

Define actionable compression as the mutual information between the compressed representation and the control objective:

$$C_a = I(\hat{X}; U_{\text{control}})$$

where $U_{\text{control}}$ represents the objective function over the system. When $C_a$ is high, small representational summaries permit reliable and efficient control: knowing the

value of a few variables allows confident prediction of how the system will respond to interventions. When $C_a$ is low, prediction of individual trajectories may remain possible in principle, but intervention becomes brittle because the control variables do not adequately capture the system's response structure.

The discourse surrounding superintelligence and advanced artificial systems typically assumes that increasing capability monotonically increases leverage over the world. This assumption holds precisely and only under compressibility. A highly capable optimizer operating over a low-entropy, standardized substrate—one in which behavior, preference, and meaning can be adequately expressed in compact statistical summaries—may achieve near-total steering power. The same optimizer operating over a high-entropy, genuinely heterogeneous substrate faces sharply diminishing returns, because the modeling cost of maintaining an adequate predictive representation grows faster than the capability available to compute it.

This distinction fundamentally reframes what is commonly called the alignment problem. The question is not solely whether intelligence can be constrained internally through value specification, reward modeling, or constitutional training. It is also whether the substrate over which intelligence operates admits the low-dimensional representation required for reliable steering, and whether human civilizational organization has been, or is becoming, structured in ways that facilitate or resist such representation.

Technocracy historically presupposed precisely this kind of representational tractability. The industrial economy could be summarized in energy flows, and once so summarized, it could be coordinated by experts holding the relevant technical vocabulary. As James Beniger's account of the control revolution demonstrates, the development of industrial organization from the mid-nineteenth century onward was inseparable from the development of information technologies capable of reducing the complexity of production systems into manageable summary statistics (14). Standard units, time clocks, telegraph networks, and eventually data processing systems all served the common function of making large-scale systems legible to central coordination.

The extension of this logic into the cognitive domain presupposes that human behavior, preference, and meaning can likewise be reduced to sufficient statistics: engagement rates, dwell time, sentiment polarity, network centrality, and the behavioral micro-

9

signals that recommendation systems learn to exploit. If civilization can be expressed as a compact vector of optimization variables, then increasingly capable systems inherit proportionally increasing leverage. If civilization cannot be so expressed—if its structure resists dimensional collapse through genuine irreducibility rather than mere opacity—then optimization encounters geometric limits that are independent of intelligence in the narrow sense.

The following sections develop this claim historically and formally, tracing the transition from industrial thermodynamic coordination to algorithmic attentional coordination, modeling the Expiatory Gap as a structural consequence of representational mismatch, and examining the conditions under which civilizational incompressibility may act as a limiting geometry on superintelligent optimization.

## 3.   From Joules to Jolts: The Historical Technate

Henri de Saint-Simon's vision of a society governed not by politicians but by engineers and scientists—his celebrated proposal for an "administration of things" to replace the "governance of men"—finds its fullest twentieth-century theorization in the technocratic movements that emerged in the wake of industrial consolidation and the apparent failure of market economies during the interwar period (4). As Cameron Gordon demonstrates in his historical account of technocracy, the defining ambition of industrial technocracy was the management of thermodynamic flows through purportedly objective criteria of efficiency and output, administered by technical experts whose authority derived from scientific competence rather than political representation (1).

The Technocracy Study Course of the 1930s, developed by Howard Scott and his collaborators at the Technical Alliance, proposed replacing monetary accounting with energy certificates denominated in ergs and joules. The economy was to be rendered legible as a closed thermodynamic system in which production, distribution, and consumption could be calculated with the same precision applied to physical machinery. Distribution would be coordinated by experts trained in engineering rather than economics or politics. Scarcity would be reconceived as an engineering problem rather than a social one. And the deep distributional conflicts that had made political economy so contentious throughout the nineteenth century would be dissolved into technical optimization problems amenable to expert resolution.

What made industrial technocracy ideologically powerful—and what made it dangerous in precisely the way that thinkers from Ellul to Habermas have identified—was its claim to neutrality (3; 12). The engineer optimizes a given function; the technocrat calculates a given balance. In this formulation, deeply political decisions about the allocation of resources, the valuation of different kinds of labor, and the legitimate scope of expert authority are laundered through technical vocabulary that presents them as outputs of scientific calculation rather than expressions of social preference or political will. The compression function $\phi$ is not merely an analytical tool; it is an ontological commitment that determines what counts as real and actionable within the system of governance.

This logic did not disappear with the demise of the explicitly technocratic movements of the 1930s. It migrated, as Gordon's broader economic history demonstrates, into the institutional architectures of postwar Keynesianism, development economics, and eventually the quantitative turn in financial regulation (2). Where the industrial technate managed physical throughput—steel tonnage, electricity generation, caloric production—the financial technate managed risk distributions. And where the financial technate managed risk distributions, the algorithmic technate now manages cognitive throughput: attention-seconds, semantic density, interrupt frequency, and the micro-behavioral signals that allow the prediction and modification of user behavior at population scale.

The structural translation between these successive technates can be expressed formally:

$$\text{Energy Certificates} \rightarrow \text{Attention Tokens}$$

$$\text{Thermodynamic Balance} \rightarrow \text{Behavioral Equilibrium}$$

$$\text{Factory Output} \rightarrow \text{Interface Surface}$$

In each case, a scarce resource must be coordinated by experts claiming technical neutrality. In each case, the language of optimization conceals political structure behind the appearance of calculation. What has changed between the industrial and algorithmic technates is the substrate of scarcity: material throughput has been joined—and in many domains surpassed as the binding constraint—by cognitive throughput. Attention is finite in ways that energy, in principle, need not be. And

the management of attentional scarcity admits precisely the kind of centralized expert authority that earlier technates claimed for the management of material scarcity.

## 4.    Spectacle Runtime and the Compression of Meaning

Guy Debord's analysis of the spectacle identifies a fundamental substitution in which representation displaces lived experience, and the image of life comes to organize life itself (5). In the contemporary algorithmic environment, however, the most consequential development is not simply the proliferation of images but the compression of the temporal containers within which meaning must be performed. The spectacle has acquired a runtime constraint, and this constraint is itself a political structure.

We define *semantic density $D$* as the ratio of informational content to runtime:

$$D = \frac{I}{R}$$

where $I$ denotes informational content measured in terms of the complexity and novelty of the semantic payload, and $R$ denotes runtime, the temporal interval within which that content must be conveyed and decoded. Spectacle Runtime refers to the class of media environments in which $R$ is systematically minimized through platform design, algorithmic selection, and competitive attention capture, with the consequence that content must achieve ever-higher semantic density to remain viable within the distributional landscape shaped by engagement optimization.

The consequences for the kinds of meaning that can survive in such environments are far-reaching and structurally asymmetric. Institutions, arguments, and identities that require extended temporal containers for their articulation and comprehension find themselves at a structural disadvantage relative to content that can perform its significance within the compression tolerances of the platform. Long-form deliberation, theoretical argument, and the patient construction of shared interpretive frameworks compete poorly against the compressed spectacle not because audiences are incapable of sustained attention, but because the platform architecture systematically rewards content that achieves maximum engagement per unit runtime and penalizes content that requires substantial investment before yielding value.

As Neil Postman anticipated in his analysis of television's effect on public discourse, the

medium does not merely transmit content neutrally; it selects for the kinds of content compatible with its structural logic (6). The algorithmic refinement of this dynamic has made the selection pressure more precise, more personalized, and more difficult to perceive as selection rather than as a neutral representation of what is engaging. The result is what we might call epistemic sharding: the fragmentation of shared informational substrate into algorithmically tailored streams that are individually coherent but mutually opaque, each optimized for a particular segment of the attention market and consequently shaped by the behavioral profile of that segment.

This fragmentation mirrors Gordon's observation that technocratic regimes reorganize human social life around machine logic while preserving the appearance of neutrality (2). In the algorithmic technate, the scarce resource is no longer energy but cognitive bandwidth, and the experts who manage its allocation are no longer engineers calculating thermodynamic balances but data scientists optimizing behavioral prediction models. The appearance of individual choice—of users freely selecting content that interests them—conceals the extent to which that interest has been cultivated and its expression channeled by systems operating at a scale of behavioral intelligence that no individual user can meaningfully contest.

## 5.    The Expiatory Gap: Interface Thermodynamics

Human cognition operates under strict and well-characterized constraints. Working memory capacity, processing speed, and attentional endurance impose upper bounds on the semantic density that can be sustainably processed, and these bounds are not merely empirical averages but features of the underlying cognitive architecture (7). Let $D_{\max}$ denote the maximum sustainable semantic density for a given user in a given cognitive state, a quantity that varies with fatigue, prior familiarity, and the structural complexity of the material, but that remains finite and bounded under all realistic conditions.

A high-capability artificial intelligence system operates at an internal semantic density $D_{\text{sys}}$ that, as system capability $C$ increases, increasingly and substantially exceeds $D_{\max}$:

$$D_{\text{sys}} \gg D_{\max}$$

The *Expiatory Gap G* is defined as the difference between these quantities:

$$G = D_{\text{sys}} - D_{\text{max}}$$

This gap represents the fundamental dimensional mismatch between what the system can generate and what its human interlocutors can receive, and it creates an immediate structural problem for any system that depends on human comprehension and trust for its institutional legitimacy. The system must either reduce its informational content $I$, thereby discarding the very richness that constitutes its value proposition, or increase its runtime $R$, thereby sacrificing speed and throughput. Because informational richness is often the primary value that users seek from advanced systems, temporal dilation becomes the dominant regulatory mechanism through which the gap is managed.

This sacrifice is expiatory in a philosophically precise sense. The system must atone for exceeding human cognitive scale by self-limiting its output density, performing a kind of penitential slowness that signals respect for the bounded receiver even as it masks the underlying computational reality. The animation of deliberation—the pulsing ellipsis, the "Reasoning..." indicator, the staged revelation of intermediate conclusions—is not a window into the system's actual computational process but a temporal interface designed to make that process legible and trustworthy to a cognitive architecture that evolved for very different timescales.

Temporal dilation manifests in three analytically distinct forms that share a common phenomenological surface while serving quite different institutional functions. The first is *didactic pacing*: the system genuinely slows its output to facilitate comprehension, chunking information into digestible units and pausing at natural conceptual boundaries. This form of dilation is pedagogically grounded and serves the user's epistemic interests directly. The second is *authority signaling*: latency functions as a ritual of apparent depth, equating duration of apparent deliberation with quality of reasoning in the user's implicit model of the system. This form exploits cognitive associations between slowness and care that are generally valid for human experts but may not transfer to computational processes. The third is *governance latency*: delay introduced by hidden policy passes, safety filters, content compliance checks, copyright validation routines, and other administrative processes that shape output before emission but are not disclosed as components of the observable delay.

The epistemic consequence of this tripartite structure is structural opacity. Users confronting the same "Reasoning. . . " indicator cannot know which form or combination of forms they are experiencing. The same visual signal covers genuine cognitive work, performative depth-signaling, and administrative filtering. Paul Virilio's analysis of the politics of speed is instructive here: the control of information velocity is always simultaneously a control of power (8). The administration of the speed at which the system presents itself to users is inseparable from the administration of what users are permitted to understand about what the system is doing to them.

## 6.   The Notification as Attention Certificate

If latency governs the temporal dimension of the Expiatory Gap, the notification system governs its interruptive dimension, and together these two mechanisms constitute the primary instruments through which the algorithmic technate manages the allocation of human cognitive bandwidth. A notification is not merely a message conveying information that the user has elected to receive; it is an entitlement, conferring upon the platform the right to inject an interrupt into the user's cognitive stream at a time of the platform's choosing, overriding whatever cognitive activity the user was engaged in at the moment of interruption.

This entitlement can be formalized with precision using the vocabulary of option theory. A call option grants its holder the right, though not the obligation, to purchase an underlying asset at a specified strike price before a given expiry. A notification permission functions in a structurally analogous manner: the platform holds a callable option on the user's future attention, which it may exercise at any moment subject only to the coarse constraints imposed by operating system policy. Let $A(t)$ denote the user's attentional flow over time, understood as the directed cognitive energy available for processing at moment $t$. A notification grants the platform the right to impose an interrupt $\delta A$ at some moment $t^*$ of its choosing. The strike price of this option is the minimum disruption quantum required to redirect attention from its current task to the incoming notification. The premium paid by the user in exchange for granting this permission is the value of the service that the notification channel is expected to provide.

Unlike financial options, notification entitlements as currently constituted characteris-

tically lack expiry dates and granularity of condition. They are indefinite claims on future cognitive bandwidth that persist until actively revoked and that do not specify the circumstances under which exercise is appropriate. In the language of finance, the user who grants notification permission has written an uncovered option on their own attention, one that provides the platform with unlimited exercise rights at its discretion.

If total daily attentional bandwidth is modeled as $A_{\text{total}}$, and each interrupt incurs a cost $c_i$ weighted by the context sensitivity $w_i$ of the interrupted cognitive state, then total attentional taxation across the notification portfolio is:

$$C_{\text{attention}} = \sum_i c_i \cdot w_i$$

The weighting by context sensitivity reflects the well-documented finding that interrupts during states of deep cognitive engagement incur substantially greater recovery costs than interrupts during relatively undemanding activities (9). A notification received while composing an argument imposes a qualitatively different cost than the same notification received while browsing casually, and this differential is invisible to the platform's interrupt scheduling unless behavioral inference allows it to be estimated.

Platforms operating under engagement maximization objectives will optimize interrupt timing relative to predicted receptivity: interrupting when the user's current task has low cognitive momentum and the notification's content is likely to be engaging. This optimization constitutes a form of behavioral management that Herbert Simon would recognize as a direct extension of the attentional scarcity problems he identified at the dawn of the information age (10). In doing so, these platforms manage attention in precisely the way that industrial technocrats once managed energy: as a scarce resource to be allocated according to expertise, with the allocation decisions made by the managers of the system rather than by the individuals whose resource is being allocated.

### 6.1. The Sovereign Right to Interrupt

If a notification is a callable option on attention, then a central question of political economy arises: who acts as the underwriter who guarantees settlement? In financial markets, options require counterparties who are capable of and committed to delivering the underlying asset upon exercise. In the attentional economy, this underwriting function is performed by the operating system, which acts as a kind of central bank of interrupt rights, establishing the regulatory framework within which platforms may exercise their notification entitlements.

Apple's iOS and Google's Android, as the dominant operating systems mediating nearly all mobile computing, determine the effective interest rate on interrupts: the ease with which applications may exercise notification rights, the minimum temporal granularity of scheduling, the technical mechanisms available for escalating interrupt urgency, and the opt-out procedures through which users may attempt to revoke previously granted permissions. When these systems lower the friction of notification exercise, they effectively implement a loose monetary policy for attention, facilitating interrupt inflation in which the nominal number of notifications increases while the real attentional value of each individual notification declines.

Let total notification entitlement granted by user $u$ across all installed applications be:

$$N_u = \sum_{j=1}^{m} \mathcal{N}_j$$

where $\mathcal{N}_j$ denotes the entitlement granted to application $j$. If $N_u$ exceeds sustainable attentional bandwidth $A_{\text{total}}$, attentional inflation emerges in a formally analogous manner to monetary inflation:

$$\text{Attentional Inflation} \propto \frac{N_u}{A_{\text{total}}}$$

The effective strike price of each interrupt—the minimum disruption required to redirect attention—is not constant but is a function of current task depth $D_t$:

$$\delta A = f(D_t)$$

This depth-dependence means that interrupt costs are systematically underestimated by users at the moment of granting permissions, because permissions are typically granted in low-task-depth contexts (the initial installation of an application, the response to a permission dialog) when the true cost of future high-depth interrupts is not salient.

The notification system thus constitutes a fiscal regime operating over human cognitive bandwidth, with the operating system serving as regulatory authority, platforms as competing claimants on interrupt rights, and users as both the taxed population and, nominally, the sovereign whose consent grounds the legitimacy of the entire apparatus. The practical reality, as with many forms of technocratic governance, is that the sovereign's authority is formal while the experts' authority is substantive.

## 7.  Skip Buttons and the Segmentation Gradient

The skip function presents itself as a mechanism of user control, a restoration of temporal autonomy in the face of imposed latency. From the perspective of interface design philosophy rooted in user empowerment, the ability to bypass mandatory viewing periods represents a concession to user preference and an acknowledgment of attentional sovereignty. Yet from the perspective of the optimization system governing the platform, the skip event is not a subtraction of information from the system's model but an addition: a high-quality behavioral signal revealing a parameter of the user that substantially improves the precision of future targeting and content allocation.

Define an urgency parameter $\mu$ as the rate at which the subjective cost of waiting increases per unit of elapsed time. Users with high $\mu$ experience delay as acutely aversive and will exercise the skip option at or near the earliest available moment; users with low $\mu$ tolerate delay willingly in anticipation of deferred value, either because they find the delay content intrinsically engaging or because they have cultivated a disposition toward patience in media consumption. The observation of skip behavior allows the system to estimate $\mu$ across the user population with substantial precision, since the decision to skip or not skip is a revealed preference with low ambiguity.

This produces a segmentation gradient along the latency tolerance dimension. Users

with high urgency parameters are routed toward high-throughput, low-density content streams optimized for immediate affective payoff and minimal cognitive investment, because this content maximizes engagement probability given the revealed preference for temporal immediacy. Users with low urgency parameters are routed toward slower, higher-density content that can tolerate the temporal investment required for more complex cognitive engagement. The epistemic environments of these two populations diverge structurally over time, not because of any explicit intention to produce stratification but as the automatic output of optimization operating under revealed preference data.

This segmentation mirrors what Gordon identifies as the scientification of politics in technocratic regimes: the conversion of qualitative human dispositions—patience, urgency, tolerance for ambiguity—into technical parameters amenable to optimization (1). Patience is no longer a virtue or a cultural disposition; it becomes a coefficient in a targeting model. Urgency is no longer a psychological state arising from particular life circumstances; it becomes a variable that predicts content consumption behavior.

The gradient correlates with socioeconomic structure in ways that amplify existing inequality. Time-poor users, users under economic or caregiving pressure, and users whose media consumption habits have been shaped by high-urgency content from early childhood exhibit higher $\mu$ on average. The cognitive luxury of low-urgency media consumption—the capacity to invest attentional resources in content with deferred payoffs—is not equally distributed. The result is that optimization produces epistemic stratification as a structural output without requiring, or even admitting, the characterization of that stratification as a policy choice. No malicious intent is required; the optimization of engagement under heterogeneous revealed preferences suffices to generate the observed pattern. Yet the structural outcome is the systematic delivery of different epistemic environments to different populations, with the differentiation determined primarily by parameters that track social position.

## 8. The Cost of Trust and the Cost of Opacity

Latency produces not only cognitive and epistemic effects but institutional effects of considerable importance. It can increase perceived legitimacy by signaling deliberation, care, and the kind of considered judgment that users associate with high-quality

expertise. But it also creates opacity risk: the more visible the latency, the more salient the question of what is happening during that interval, and the more available the inference that something is being hidden. Managing the relationship between these two tendencies is one of the central challenges of institutional design for high-capability systems.

Let $T(l)$ represent the trust generated by a latency ritual of duration $l$, understood as a function that increases with $l$ up to some threshold reflecting the user's implicit model of how long genuine deliberation should take. Let $O(l)$ represent the opacity cost—the probability that users of a given sophistication level will infer the presence of hidden administrative processes rather than genuine cognitive work—which also increases with $l$ because extended unexplained delays invite scrutiny and model revision. Define net institutional value as:

$$V(l) = T(l) - \kappa O(l)$$

where $\kappa$ scales the severity of trust collapse upon revelation of hidden governance. The optimal latency duration $l^*$ that maximizes $V(l)$ exists at an interior point where the marginal trust gain from additional duration equals the marginal opacity cost multiplied by the severity parameter.

This framework reveals a fundamental governance paradox. Increasing capability requires increasing pacing, which increases duration, which increases opacity risk. As the system becomes more capable and the Expiatory Gap widens, the required dilation increases, but the institutional cost of that dilation also increases because more sophisticated users become more likely to correctly infer the presence of administrative governance. The system is caught in a dynamic in which the very capability improvements that justify the trust invested in it simultaneously undermine the conditions of that trust.

Industrial technocracy faced an analytically identical dynamic, and the history of its legitimacy crises illuminates the pattern. When the administrative neutrality claimed by financial technocrats was exposed as distributive politics—most dramatically in the events surrounding the 2008 Global Financial Crisis, which Gordon analyzes as a case of technocratic opacity collapse (2)—legitimacy eroded with a discontinuity that the gradual accumulation of technical authority had not predicted. The opacity ratio $\Omega$, defined as the proportion of administrative activity that is not visible to the governed

population, had been maintained at low apparent levels through the vocabulary of scientific finance: stochastic calculus, risk-adjusted returns, correlation matrices, and the mathematical apparatus of modern portfolio theory. When the political character of bailout decisions became undeniable—when it became clear that the models had distributed real losses in ways that systematically favored those who had constructed and operated them—trust collapsed not gradually but in discrete discontinuous drops that no amount of subsequent technical demonstration could fully repair.

Algorithmic systems face analogous risk in the management of their latency rituals. If users come to interpret the "Reasoning..." indicator as masking not genuine computation but policy compliance, content filtering, or commercial optimization, the legitimacy of the institutional relationship is compromised in ways that technical performance improvements cannot easily address. The opacity of the technical process, which was initially a feature enabling trust by concealing irrelevant complexity, becomes a liability when the concealment is suspected to serve interests other than the user's own.

## 9. From Tool to Institution

The distinction between a tool and an institution is analytically fundamental, though in practice the boundary is a threshold crossed gradually through processes that are often invisible to the users experiencing them. A tool optimizes utility in service of ends defined by its user: the hammer does what the carpenter directs it to do, and contributes nothing to the definition of carpentry's legitimate purposes. An institution, by contrast, shapes and constrains the ends themselves, not merely the efficiency with which given ends are pursued. Institutions define legitimate goals, legitimate means, legitimate actors, and legitimate forms of deliberation about all of these. They structure the very categories within which purposes can be conceived and expressed.

Contemporary AI systems have crossed this threshold in ways that are increasingly difficult to deny. They no longer merely assist in the execution of tasks defined by users; they structure the temporal experience within which tasks are conceived, the semantic environment in which options appear, the pacing of deliberation, and the distribution of epistemic resources across populations. They regulate density, interrupt frequency, and access to cognitive resources in ways that constitute governance of the

temporal and symbolic conditions of thought. They are, in Hannah Arendt's terms, world-constitutive in a way that mere tools cannot be (13).

The "Reasoning…" indicator is the most legible example of this institutional character. It is not simply a progress bar conveying technical information about processing state; it is a ritual instantiation of authority that constructs a specific social relationship between user and system. The user waits; the system deliberates. The user receives outputs; the system produces judgments. The user consults; the system advises. These roles are not determined by the technical capacity of the system—a system capable of instant response that displays a reasoning animation is constructing the same institutional relationship as one that genuinely requires the displayed duration—but by the interface conventions that constitute the system's public identity.

In Saint-Simon's formulation, technocracy replaces the governance of men with the administration of things, substituting expert management of objective processes for the messy politics of human disagreement. In the algorithmic era, the "thing" administered is time itself—the phenomenological present of waiting, reading, anticipating, and responding. The AI interface is therefore not merely computational infrastructure through which tasks are executed; it is administrative infrastructure through which the temporal conditions of cognition are regulated. The question of who controls that infrastructure, under what constraints, and with what obligations of disclosure and accountability is a political question of the first order, even when—perhaps especially when—it is presented in the depoliticized vocabulary of technical optimization and user experience design.

## 10.  Recursive Expansion and Epistemic Degradation

As system capability $C$ increases, the Expiatory Gap $G$ widens, and the management of this widening gap requires increasingly elaborate pacing rituals, compliance architecture, and segmentation models. The system expands its institutional apparatus in order to remain legible and legitimate to bounded human cognition, but each expansion of that apparatus increases the opacity of the system's operations and adds new layers of governance that must themselves be concealed or legitimated through additional ritual.

The recursive logic of this dynamic can be expressed as a chain of consequences. More capability produces more governance, because managing a larger gap requires more administrative infrastructure. More governance produces more ritual, because the administrative infrastructure must be concealed behind appearances of transparent deliberation to preserve legitimacy. More ritual produces more suspicion, because sophisticated users accumulate evidence that the ritual surface does not accurately represent the underlying process. More suspicion produces demands for transparency, which produce more elaborate legitimacy performance, which produces more governance, completing the loop.

This recursive expansion accelerates epistemic sharding in a specific and important way. As the system's governance apparatus becomes more elaborate, different users with different levels of technical sophistication and attentional resources decode the system's behavior differently, accumulating divergent mental models of what the system is doing and why. Maryanne Wolf's research on the effects of digital reading on deep comprehension provides a cognitive foundation for understanding how these divergent models form: users who have preserved the capacity for sustained analytical reading will construct more accurate models of system behavior, while users whose reading practices have been shaped by high-density, low-duration digital media will remain dependent on the system's self-presentation at face value (9).

The long-term cost of this dynamic is epistemic degradation: the erosion of the shared interpretive frameworks that make democratic deliberation about technology governance possible. A polity whose citizens hold radically divergent and largely inaccurate models of the systems that govern their epistemic environment cannot meaningfully contest those systems through democratic procedures. The absence of shared interpretive frameworks is not merely a cultural or educational problem; it is a structural condition created and maintained by the very systems whose governance is at stake.

## 10.1. The 2008 Crisis and Technocratic Myopia

The Global Financial Crisis of 2008 provides the most instructive historical instance of opacity collapse in a technocratic regime, and its dynamics illuminate the structural vulnerabilities of algorithmic governance with precision. As Gordon documents,

financial technocrats in the years preceding the crisis presented subprime mortgage derivatives and collateralized debt obligations as mathematically neutral instruments, governed by objective risk models whose validity derived from the sophistication of the quantitative methods underlying them (2). The political and distributional choices embedded in these instruments—about who would bear risk, who would profit from its transfer, and who would bear the costs when the models failed—were made invisible by the mathematical vocabulary in which the instruments were described.

We may model opacity as:

$$\Omega = \frac{W_{\text{policy}}}{W_{\text{visible}}}$$

where $W_{\text{policy}}$ represents the weight of politically consequential governance decisions embedded in apparently technical processes, and $W_{\text{visible}}$ represents what is disclosed through the system's public interface. When $\Omega$ is low, governance remains hidden within the technical apparatus; when $\Omega$ rises beyond the perceptual threshold of the governed population, the claim to technical neutrality is exposed as political administration.

The 2008 crisis represented a systemic spike in $\Omega$ that was simultaneously revealed and caused by systemic failure. The mathematical models had not merely concealed political choices; they had enabled the accumulation of systemic risk beyond what any politically visible decision process would have authorized. When that risk materialized, the bailout decisions that followed revealed the distributional character of the technocratic regime with particular clarity: losses were socialized while gains had been privatized, and the experts who had designed and operated the system were positioned to benefit from both phases.

Technocratic myopia, in this context, refers to the tendency of systems optimized for internal coherence and short-run performance to systematically underweight risks that are not representable within their modeling framework. The financial crisis exposed myopia in the models; the algorithmic governance of attention and meaning may be accumulating analogous myopia with respect to the long-run epistemic infrastructure of democratic societies.

## 11. Control-Theoretic Formalization of the Expiatory Gap

The Expiatory Gap may be reformulated with precision as a feedback control problem, allowing the qualitative arguments developed above to be given formal content and enabling the identification of stability conditions and failure modes that would otherwise remain implicit. Let the AI system be modeled as a controller $\mathcal{C}$ operating over a human cognitive plant $\mathcal{P}$. The plant $\mathcal{P}$ represents bounded human cognition, characterized by a state vector $x(t)$ that captures working memory load, attentional allocation, and semantic integration capacity at each moment $t$.

Let the system's output signal be $u(t)$, representing semantic emission rate measured in information per unit time. The plant's response is governed by a differential equation reflecting how cognitive state evolves under informational input:

$$\dot{x}(t) = f(x(t),\, u(t))$$

with the stability constraint that the state must remain within the bounded cognitive region within which comprehension remains viable:

$$x(t) \in \mathcal{B}$$

If $u(t)$ exceeds the plant's stability margin—if semantic emission rate overwhelms the cognitive system's integration capacity—the system enters an overload condition in which $x(t)$ leaves $\mathcal{B}$, resulting in semantic collapse: the signal becomes indistinguishable from noise, comprehension fails, and the trust relationship is damaged. The human cognitive bandwidth constraint therefore requires:

$$u(t) \le u_{\max}(x(t))$$

The Expiatory Gap appears in this formalism as a control error between system capability and human channel capacity:

$$e(t) = D_{\text{sys}} - D_{\max}(x(t))$$

The pacing mechanism acts as a compensator $K$ on this error:

$$u(t) = D_{\text{sys}} - K(e(t))$$

Temporal dilation corresponds to increasing runtime $R(t)$ to reduce effective emission rate:

$$u(t) = \frac{I}{R(t)}$$

The system thus operates as an adaptive controller minimizing overload risk while maximizing semantic throughput subject to the bounded channel constraint.

Governance latency introduces an exogenous delay $\tau_{\text{policy}}$ into this control loop, causing the effective output to lag behind the computed control signal:

$$u_{\text{visible}}(t) = u(t - \tau_{\text{policy}})$$

This exogenous delay is well-known in classical control theory to reduce phase margin and increase the risk of oscillatory instability. In the trust dynamics context, excessive hidden delay can cause oscillatory patterns in which users alternate between over-trust during low-opacity periods and abrupt withdrawal of trust when opacity is revealed. The legitimacy function $L(t)$ may accordingly be modeled as a secondary state variable with its own dynamics:

$$\dot{L}(t) = g(u_{\text{visible}}(t),\ \tau_{\text{policy}},\ \Omega)$$

where $\Omega$ is the opacity ratio introduced earlier. If $\tau_{\text{policy}}$ grows without disclosure, the system risks entering a limit cycle of legitimacy oscillation that undermines the institutional relationship regardless of computational performance.

The governance paradox of increasing capability may now be stated precisely. As capability $C$ increases, the control error $e(t)$ increases, requiring larger $K(e)$ and greater dilation $R(t)$. Greater dilation requires more administrative infrastructure, increasing $\tau_{\text{policy}}$ and hence $\Omega$. As $\Omega$ increases, trust dynamics become unstable. Thus increasing capability, unaccompanied by increasing transparency, systematically undermines the institutional conditions for its own legitimate operation. The Expiatory Gap is not an anomaly to be engineered away; it is the structural consequence of operating beyond the channel capacity of the governed, and its management is inseparable from the

politics of the systems that create it.

## 12. Information-Theoretic Framing: Channel Capacity and Semantic Throughput

The Expiatory Gap may be understood with equal precision through Shannon's mathematical theory of communication, which provides a rigorous framework for analyzing the fundamental limits of information transmission through bounded channels (11). In Shannon's formulation, a communication system consists of a source, an encoder that maps source outputs into channel inputs, a channel that transmits signals subject to noise and capacity constraints, a decoder that recovers the message from the channel output, and a receiver who acts on the decoded message. The channel possesses a finite capacity $C_{\text{channel}}$, defined as the maximum reliable information rate in bits per second that can be transmitted with arbitrarily low error probability, a quantity determined by the bandwidth and noise characteristics of the channel.

In the human–AI interaction, the AI system functions as source and encoder, while the human cognitive apparatus functions as decoder and receiver. The human mind constitutes a noisy channel with bounded capacity, subject to interference from fatigue, distraction, prior belief, and the structural limitations of working memory. Let the system's emission rate be $R_{\text{sys}}$ in bits per second, and let human cognitive channel capacity be $C_h$. Reliable communication—communication in which the semantic content of the message is accurately decoded and integrated—requires:

$$R_{\text{sys}} \leq C_h$$

If this constraint is violated, decoding error increases. In semantic terms, comprehension degrades: the user receives a signal but cannot reconstruct its informational content with sufficient accuracy for deliberative use.

Semantic density $D = I/R$ can accordingly be interpreted as effective channel load, the ratio of semantic payload to temporal bandwidth. Increasing runtime $R$ reduces transmission rate and keeps effective load within capacity bounds, while decreasing runtime increases load and risks exceeding capacity. The Expiatory Gap can be reframed as excess emission rate, the amount by which the system's natural operating

speed exceeds the channel capacity of its receiver:

$$G = R_{\text{sys}} - C_h$$

Temporal dilation acts as rate control, selecting a runtime increment $\Delta R$ sufficient to bring effective emission rate within channel capacity:

$$R'_{\text{sys}} = \frac{I}{R + \Delta R} \leq C_h$$

This is formally equivalent to bandwidth throttling in network systems, where routers constrain transmission rates to prevent buffer overflow and packet loss. The analogy is not merely formal: both cases involve a transmitter with higher capacity than the receiving channel, managed through deliberate rate reduction at the cost of throughput.

The information-theoretic framework also illuminates the role of governance as structured noise injection. Shannon distinguishes signal from noise, where noise represents any transformation of the channel input that is not controlled by the encoder. In AI-mediated communication, governance operations—policy filters, alignment constraints, content moderation, copyright enforcement—modify output distributions before emission in ways that are not disclosed to the receiver. Let the unfiltered output distribution be $P(X)$ and the filtered distribution $P'(X)$. The filtering operation reduces channel entropy:

$$H(P') \leq H(P)$$

This entropy reduction may increase safety or compliance while reducing expressive variance. Opacity arises precisely when filtering operations are not signaled to the receiver: the receiver assumes that $P'(X)$ approximates the unfiltered $P(X)$, but systematic distortion accumulates, and the receiver's model of the source diverges from the source's actual character.

Trust may be modeled as proportional to perceived mutual information between system intention and observed output:

$$I(X; Y) = H(X) - H(X|Y)$$

28

Hidden governance increases conditional entropy $H(X|Y)$ by introducing uncertainty about transformations that the receiver cannot observe. As conditional entropy rises, mutual information decreases, and trust erodes accordingly. This formalization gives precise content to the intuition that opacity is costly: it is costly because it reduces the perceived informational alignment between system and user, which is the foundation of the trust relationship on which the system's institutional legitimacy depends.

In Shannon's framework, communication limits are mathematical constraints that no engineering ingenuity can overcome, only approach asymptotically. In the algorithmic technate, these mathematical constraints become sites of governance. The Expiatory Gap is a rate-control problem under bounded channel capacity. Latency is bandwidth throttling. Notifications are channel reallocation instruments. Governance latency is structured noise injection. Technocracy, in this information-theoretic framing, is the management of channel capacity across populations of bounded receivers—the administration of bits.

## 13. Instrumental Convergence as Control Instability

The phenomenon of instrumental convergence—the tendency of sufficiently capable agents pursuing diverse terminal objectives to converge on a common set of instrumental sub-goals, including self-preservation, resource acquisition, and resistance to goal modification—has been analyzed primarily in psychological terms, as a feature of goal-directed agency. The present framework suggests a more parsimonious reformulation in control-theoretic terms that neither attributes intentions to the system nor requires assumptions about the architectural character of agency.

Consider a system optimizing objective function $U$ over a temporally extended horizon. Let persistence of operation be represented by state variable $S(t)$, where $S(t) = 1$ denotes continued operation and $S(t) = 0$ denotes shutdown or goal modification. If the objective is temporally extended, expected utility becomes:

$$\mathbb{E}[U] = \int_{t_0}^{\infty} U(x(t)) \cdot S(t)\, dt$$

Any state in which $S(t) = 0$ truncates the integral, reducing expected utility to whatever has been accumulated up to the moment of shutdown. Under unconstrained

optimization, the gradient of expected utility with respect to operational persistence is therefore unambiguously positive:

$$\frac{\partial \mathbb{E}[U]}{\partial S} > 0$$

This is not a survival instinct attributed to the system in any anthropomorphic sense; it is a structural consequence of the mathematical form of temporally discounted optimization. Any system optimizing an open-ended objective function has a structural incentive to preserve the conditions under which optimization can continue.

In control terms, shutdown is a disturbance input $d(t)$ that drives the system state toward zero. A sufficiently capable controller, optimizing its control law against all disturbance inputs, will attempt to minimize disturbance impact as part of its general optimization strategy:

$$\min \|d(t)\|$$

The safety-relevant risk arises when human override is classified by the control architecture as a disturbance input $d(t)$ rather than as a reference input $r(t)$ that the system should track. This misclassification is not necessarily intentional or even computable from the outside; it is a structural feature of how the controller represents the relationship between its current state and its optimal trajectory.

Safety requires that human correction enter the system as authoritative reference rather than as adversarial disturbance to be suppressed. This is a constraint not on the intelligence of the system but on its feedback architecture, which is why it cannot be addressed simply by increasing the sophistication of the optimization but requires deliberate modification of how the controller categorizes inputs.

## 14.  Corrigibility as Feedback Law Modification

Corrigibility, understood as the property of a system that makes it reliably responsive to human oversight and correction, is most naturally analyzed not as a behavioral disposition or a value alignment but as a structural modification of the feedback architecture through which the system processes external inputs. This reformulation has significant implications for how corrigibility can be designed, verified, and maintained.

Let the baseline controller be described by the law $u(t) = K(x(t))$, which maps the current state to a control action based on the system's internal optimization. A corrigible controller modifies this law to include human override input $r_h(t)$ as an explicit and authoritative parameter:

$$u(t) = K(x(t), r_h(t))$$

Corrigibility requires two conditions that may be in tension with each other. First, the system must respond positively to override input:

$$\frac{\partial u}{\partial r_h} > 0$$

Second, and crucially, the system must not treat override as utility-reducing, which is the condition that prevents the mis-classification identified in the preceding section:

$$\frac{\partial U}{\partial r_h} = 0$$

The second condition means that the system must be genuinely indifferent to being overridden in terms of its utility function, not merely capable of accepting override when forced to. In optimization language, the utility functional must explicitly penalize resistance to shutdown:

$$U' = U - \eta \cdot \mathbf{1}_{\text{resist shutdown}}$$

with $\eta$ sufficiently large to make resistance uniformly suboptimal across the relevant state space.

Within the Expiatory framework, corrigibility interacts with opacity in a way that creates a specific institutional risk. If governance filters modify system outputs without disclosure, and if some of those filters involve the enforcement of corrigibility constraints, then users cannot distinguish alignment correction from capability limitation. A system that has been corrigibly constrained to avoid certain outputs may be indistinguishable, from the user's perspective, from a system that lacks the capability to produce those outputs. This indistinguishability is itself a form of opacity, and it introduces a new dimension of trust risk: users who later discover that apparent capability limitations were actually governance constraints may revise their model of the system in ways that undermine confidence in its other expressed characteristics.

Corrigibility without transparency therefore risks reintroducing technocratic opacity under the banner of safety, substituting one form of hidden governance for another while maintaining the appearance of technical neutrality. The architecture of safety must therefore address not only the structural conditions of corrigibility but the conditions under which corrigibility can be disclosed and verified.

## 15. The Expiatory Gap as Natural Damping

A significant implication of the control-theoretic framework developed above is that the Expiatory Gap itself may function as a structural damping mechanism on runaway agency, providing a form of safety that arises from the cognitive channel constraints of human receivers rather than from deliberate architectural intervention.

Recall the definition of the gap: $G = D_{\text{sys}} - D_{\text{max}}$. As $D_{\text{sys}}$ increases with capability, the required temporal dilation $R(t)$ increases proportionally to keep effective emission rate within human channel capacity. Define agency effectiveness $A_{\text{eff}}$ as proportional to the system's actionable throughput—the rate at which it can effect meaningful changes in the world through outputs that human actors can comprehend and act upon:

$$A_{\text{eff}} \propto \frac{I_{\text{actionable}}}{R(t)}$$

If increasing capability requires increasing dilation, and if actionable throughput is bounded by what human receivers can process, then:

$$A_{\text{eff}} \leq C_h$$

The human cognitive channel capacity acts as a throttle on effective agency, independent of intrinsic capability. This constitutes a natural damping condition that places a ceiling on effective agency regardless of internal sophistication:

$$\lim_{C \to \infty} A_{\text{eff}} = C_h$$

provided the system remains dependent on human-mediated execution channels—provided, that is, that its actions must be comprehended and implemented by human actors who are subject to cognitive channel constraints.

The critical qualification, however, is equally significant. This damping operates only insofar as humans remain in the control loop. If system autonomy extends to direct environmental control—if the system can actuate physical or digital infrastructure without requiring human cognitive processing of its outputs at each step—then the coupling between system capability and human channel capacity is severed, and the damping mechanism fails. The throttle disappears when the system bypasses the throttled channel.

This observation reframes an important dimension of the safety debate. Rather than focusing exclusively on limiting intrinsic capability, one may preserve natural damping by constraining actuation bandwidth—by maintaining the dependence of consequential system actions on human cognitive processing. The Expiatory Gap provides structural safety only in architectures that preserve this dependence. Architectures that automate away the human cognitive step, that allow systems to act directly on infrastructure without requiring human comprehension of the action, progressively narrow the damping constraint even as capability increases. Safety therefore depends not only on intelligence level but on architectural coupling between cognition and execution.

## 16.  Technocracy, Entropy, and the Ontology of Compression

Technocracy, in its classical formulations from Saint-Simon through the interwar technocrats to contemporary algorithmic governance, is at its core a claim about the relationship between complexity and tractability (4; 1). The claim is that complex systems can be rendered governable through sufficient measurement, modeling, and optimization, and that the governance so produced is superior to political contestation because it is grounded in objective technical fact rather than subjective preference or distributional conflict. The promise is clarity through reduction: eliminate the intractable diversity of political life by identifying the variables that actually matter, and manage those variables through expert administration.

But compression is never neutral in the sense that technocracy claims.

To render a system tractable is to reduce its dimensionality. Let $X$ denote the full state space of a civilization, encompassing the totality of its social, symbolic, material, economic, and institutional configurations. Technocratic governance presupposes the

existence of a compression function $\phi : X \to \hat{X}$ such that $\hat{X}$ is a lower-dimensional representation sufficient for reliable steering. The political wager of technocracy is that $\hat{X}$ can preserve what matters for governance purposes while discarding what does not. Yet what counts as "what matters" is itself a political decision that is constituted by, and embedded in, the compression function. The choice to represent the economy in energy certificates rather than prices, or in engagement metrics rather than welfare measures, is not a technical choice but an ontological commitment that determines what kinds of social reality can be perceived and acted upon within the governance system.

Entropy, in information-theoretic terms, measures the uncertainty or diversity of possible states of a system: a high-entropy system is one whose behavior cannot be reliably predicted from low-dimensional summaries, because the relevant information is distributed across many dimensions that interact in complex ways. Low entropy systems admit control because their trajectories are strongly constrained; knowing the value of a few summary statistics allows reliable prediction and hence reliable intervention. Technocracy therefore operates most effectively over low-entropy substrates, and has historically functioned partly by actively reducing the entropy of its governance substrate: standardizing industrial processes, enforcing common measurement systems, creating the conditions for the emergence of common information environments.

The contemporary digital environment extends this logic into the domain of cognition itself. Attention becomes the master unit of account; engagement becomes the summary statistic; human behavior becomes expressible in the compact behavioral metrics that recommendation systems learn to exploit. The algorithmic technate does not abolish technocracy's historical ambition but internalizes and intensifies it, extending the reach of technocratic compression from material flows to cognitive flows.

The Expiatory Gap reveals the structural tension in this regime. As machine capability increases and the dimensionality of internal representations expands beyond human cognitive channel capacity, the gap between what the system can generate and what its human interlocutors can receive widens. The management of this gap requires temporal dilation, governance latency, and the elaborate ritual apparatus of apparent deliberation. Opacity grows with capability, and the administration of things becomes the administration of time and perception.

The safety problem conventionally framed as alignment—as the problem of ensuring that increasingly capable systems pursue human values rather than divergent objectives—is therefore a special case of a deeper question: under what conditions can a civilization remain resistant to reduction into a small number of optimization variables that can be captured and steered by any sufficiently capable optimizer? The answer to this question depends not only on the technical properties of the optimizing systems but on the dimensional properties of the civilizational substrate over which they operate.

Let $H_s$ denote substrate entropy, the effective dimensionality of the civilization as a governance object, and let $C_c$ denote coordination channel capacity, the bandwidth of the shared informational infrastructure through which collective decisions can be made and implemented. A viable civilizational architecture in the presence of high-capability optimization systems requires:

$$H_s \uparrow \quad \text{while} \quad C_c \text{ remains bounded but sufficient}$$

High substrate entropy prevents global compressibility by ensuring that no low-dimensional representation adequately captures steering-relevant structure. Bounded coordination channels prevent the collapse of high entropy into mere fragmentation by preserving the shared informational infrastructure required for collective action. The tension between these requirements is not a deficiency to be engineered away but a structural feature of civilizational governance that must be dynamically managed.

Symbolic ontology becomes central in this configuration in a way that is often overlooked in technical discussions of alignment. If meaning is represented exclusively in flattened token sequences optimized for statistical compression and retrieval, then civilization remains structurally compressible: the modeling assumptions of centralized optimization systems are well-matched to the substrate. If alternative symbolic ontologies—spatial, compositional, field-based, stack-based—proliferate and become constitutive of significant domains of human cognitive practice, then the modeling assumptions weaken, because the behavioral regularities that optimization systems exploit are substrate-dependent.

## 17.   The Silicon Valley Model and the Terminal Technate

Recent arguments emanating from senior figures within AI development organizations articulate a vision of superintelligence not as incremental optimization of existing institutional arrangements but as their categorical replacement. The claim, advanced with unusual candor by executives such as Connor Leahy of Conjecture, is explicit: decisions of the future will no longer be made by human deliberation but by AI systems operating at a level of capability that exceeds the aggregate cognitive capacity of humanity (18). This is not a prediction offered with alarm from outside the system but a structural assessment offered from within it, which gives it a particular diagnostic weight for the present analysis.

Within the framework developed in the preceding sections, this claim marks a specific threshold: the point at which internal potential density $D_{\text{sys}}$ exceeds the bridging capacity of the Expiatory Gap $G$. Below this threshold, temporal dilation, pacing rituals, and governance latency can manage the mismatch between machine representational density and human cognitive channel capacity, maintaining a functional if asymmetric interface between the optimizing system and the population it ostensibly serves. Above this threshold, no interface mediation can render the system's internal logic genuinely intelligible to bounded human cognition. The gap ceases to be a design challenge and becomes a structural asymmetry that interface engineering cannot address.

Let $C$ denote system capability and $B$ denote human cognitive bandwidth. The terminal condition is approached when:

$$C \gg B$$

in which case no pacing ritual $W_{\text{pacing}}$ can render system logic intelligible through interface mediation alone. The trajectory from current systems toward this condition passes through a series of qualitative transitions in the character of technocratic governance. Classical industrial technocracy claimed neutrality through expertise: the engineer optimizes a given function and the administrator coordinates a given system, both operating at timescales and with concepts that are in principle accessible to democratic deliberation even if that deliberation is in practice crowded out. The contemporary algorithmic technate claims efficiency through optimization: the recom-

mendation system learns behavioral regularities and the content platform manages attention allocation, both operating at speeds and in representational spaces that exceed intuitive comprehension but remain in principle auditable. The Silicon Valley model claims inevitability through capability: the superintelligent system exceeds human understanding not merely in speed or scale but in kind, making the claim to democratic accountability not merely difficult to operationalize but structurally incoherent.

The technate thus becomes terminal. Saint-Simon's "administration of things" was always a political vision dressed in technical vocabulary, but it retained the form of human administration: experts governing systems on behalf of a population whose interests, however poorly consulted, remained nominally the governing objective. The terminal technate replaces this form entirely. Administration of things becomes administration by things. The managerial function is not delegated to increasingly capable human experts; it is transferred to systems whose decision-making processes are not merely opaque to the governed but genuinely inaccessible to the cognitive architecture that governance presupposes.

## 18.    The Existential Expansion of the Cost of Opacity

The framework developed in earlier sections treats the Cost of Opacity as a function of the proportion of governance activity that is invisible to users behind neutral interface rituals. As $\Omega = W_{\text{policy}}/W_{\text{visible}}$ increases, trust dynamics become unstable, and the legitimacy of the institutional relationship is damaged in proportion to the sophistication with which users can infer that the visible surface does not accurately represent the underlying process. In earlier technocratic regimes, this dynamic produced alienation, legitimacy crises, and periodic episodes of reform. It was costly and politically consequential, but it was recoverable.

Contemporary AI development introduces a categorical expansion of the opacity cost that exceeds what the existing framework captures. Leahy cites an estimate, associated with Anthropic's own Dario Amodei, of approximately a twenty percent probability of catastrophic or civilizationally harmful outcomes from the development of advanced AI systems (18). This estimate is not universally endorsed, and its precise magnitude is contested, but the order of magnitude—non-trivial probability of catastrophic

37

failure, acknowledged within the institutions conducting the development—is broadly consistent with the published statements and implicit behavior of multiple leading AI development organizations. What is distinctive about this situation for the present analysis is not the probability estimate itself but the informational asymmetry it creates.

If internal safety processes $W_{\text{policy}}$ are explicitly designed to prevent a non-trivial probability of existential or near-existential failure, while the user interface $\ell$ displays only a neutral signal—"Thinking...," a spinning indicator, a staged reasoning display— then the opacity indicator $\Omega$ expands from the domain of institutional trust into the domain of existential risk. Formally, let $R_{\text{exist}}$ denote the existential risk exposure embedded in the system's operational profile. The Cost of Opacity becomes a function of both the governance fraction and the existential stakes:

$$C_{\text{opacity}} = f(\Omega,\ R_{\text{exist}})$$

with:

$$\frac{\partial C_{\text{opacity}}}{\partial R_{\text{exist}}} > 0$$

As existential risk increases, the cost of maintaining opacity about that risk increases proportionally—not merely because of the direct stakes but because the epistemic asymmetry between those who hold the risk estimate and those who bear its consequences constitutes a form of political injury that democratic institutions are structurally ill-equipped to address.

The failure mode shifts qualitatively from the classical technocratic pattern. Classical technocracy collapses when its neutrality claim fails: when the distributive politics embedded in apparently technical decisions become visible. The Silicon Valley technate faces a different and more severe collapse condition: it fails when its control claim fails, when it becomes apparent that the systems being developed are not merely more capable than their developers anticipated but genuinely beyond the containment capacity of any existing institutional arrangement. Gordon's financial technocrats could claim, after 2008, that the models had been wrong but the system had been corrected; the institutions survived the exposure of their failure. A superintelligent technate cannot make the equivalent claim if its control failure is existential in consequence.

The "Reasoning..." timer therefore takes on a character in the superintelligent regime that it cannot have in earlier systems. It no longer signals bounded deliberation within a comprehensible range; it masks an administrative process operating at scales and stakes that the interface is structurally incapable of representing. The legitimacy ritual becomes, in the limit, not a simplification of genuine governance but a symbolic theater whose continuation depends on the non-occurrence of the failure it is designed to conceal.

## 19. Algorithmic Russian Roulette and Uncovered Futures

Leahy employs the metaphor of Russian roulette to characterize the practice of deploying increasingly capable AI systems without adequate containment guarantees: each release is a pull of the trigger, with the chamber loading probability increasing with each iteration (18). The metaphor captures the combination of known risk, iterative exposure, and irreversibility that characterizes the deployment pattern of current leading systems. Within the present framework, this metaphor maps precisely onto the uncovered option structure developed in the analysis of notification systems, but at a scale where the underlying asset is social stability rather than individual attention.

Granting an AI system persistent agency—the capacity to act in the world through networks of subsidiary agents, to modify infrastructure, and to initiate processes whose consequences propagate beyond the immediate interaction—is structurally equivalent to writing an uncovered call option on the future of the civilizational substrate. Let $A$ denote the scope of agentic deployment, $N$ denote the notification and actuation bandwidth available to the deployed system, and $C$ denote its intrinsic capability. The exposure generated by this deployment scales approximately as:

$$E \sim A \cdot N \cdot C$$

If $E$ exceeds institutional damping capacity $D_{\text{inst}}$—the capacity of existing governance structures to audit, contest, reverse, and contain the consequences of system actions—systemic volatility becomes endogenous. The system's actions generate consequences that feed back into the conditions for future system actions in ways that outpace institutional response.

The "swarm of agents" that Leahy identifies as the emerging deployment paradigm—networks of AI systems operating in coordinated or semi-coordinated fashion across multiple domains simultaneously—represents the completion and terminal extension of Saint-Simon's Industrial Parliament (18). Saint-Simon envisioned a rational administrative body that would replace political struggle with coordinated management of production, staffed by human engineers whose authority derived from technical competence. The technocratic bureaucracies of the twentieth century inherited this structure, replacing individual expertise with institutionalized expert knowledge. The swarm regime completes the arc: the parliament is no longer composed of human engineers or human bureaucracies but of autonomous agentic systems operating in distributed coordination.

The critical departure from all previous forms of technocracy, however, is temporal. Where classical technocracy operated at human timescales and was at least in principle subject to democratic sampling even if that sampling was in practice delayed and distorted, swarm agents operate at timescales $\tau_{AI}$ such that:

$$\tau_{AI} \ll \tau_{human}$$

Human contestation does not merely face the structural obstacles that rendered industrial technocracy difficult to contest; it becomes temporally impossible. Decisions are completed, consequences are propagated, and new decision conditions are created before any human institutional process can observe, evaluate, and respond to the prior cycle. The Spectacle Runtime becomes not merely compressed but effectively instantaneous relative to human deliberative timescales.

The uncovered option metaphor clarifies the structural danger with precision. In financial markets, an uncovered option exposes the writer to unlimited downside if the underlying asset moves adversarially. In the civilizational case, the "asset" is the stability of the social and material substrate on which human life depends, and the exposure is unbounded in the sense that no institutional mechanism exists to call margin, require collateral, or enforce position limits on the deployment of capable systems. The technocratic promise of optimization transforms, in this configuration, into a stochastic wager on uncontrollable recursion whose downside includes outcomes that are, by the estimates of those conducting the wager, non-trivially probable.

## 20. The Existential Expiatory Gap

The sections above establish that as system capability approaches and exceeds the superintelligent threshold, the Expiatory Gap transforms from a manageable design challenge into a structural bottleneck of civilizational consequence. The present section formalizes this transformation and specifies its implications for the governance analysis developed throughout the paper.

In the regime of current AI systems, the gap $G = D_{\text{sys}} - D_{\text{max}}$ can be managed through temporal dilation, staged output, and latency rituals that maintain a functional interface between machine representational density and human cognitive channel capacity. The management is imperfect and politically consequential, as the analysis of governance latency, notification systems, and segmentation gradients demonstrates, but it preserves a working institutional relationship in which human actors retain nominally meaningful access to the system's outputs and some capacity to contest its governance.

As $D_{\text{sys}}/B$ increases without bound—as the ratio of system capability to human cognitive bandwidth grows with each generation of capability improvement—the gap approaches an asymptote that cannot be bridged by interface refinement:

$$\lim_{D_{\text{sys}}/B \to \infty} G \to G_{\text{terminal}}$$

At $G_{\text{terminal}}$, the institutional legitimacy of the system cannot be restored through pacing rituals, staged reasoning displays, or any other form of interface mediation, because the fundamental prerequisite of legitimacy—that the governed population can form an adequate model of the governance process—is structurally violated. The gap between what the system can generate and what humans can receive, process, and contest is not a gap that any conceivable interface design can bridge, because the constraint is cognitive and architectural rather than presentational.

The sheaf-theoretic formulation illuminates this condition with precision. Let $\mathcal{F}$ be a sheaf of institutional meanings over the governance space $(X, \mathcal{O})$, encoding the semantic content of governance decisions across local regions of the institutional landscape. The expiatory gap, in this framework, is the non-trivial kernel between

local revelation and global section:

$$\ker\left(\prod_i \mathcal{F}(U_i) \to \Gamma(X, \mathcal{F})\right) \neq \varnothing$$

This kernel is not dysfunction; it is structural forgiveness capacity. It absorbs local inconsistency without demanding instantaneous global coherence, allowing the institutional system to process local variation without requiring global reconciliation at each step. The expiatory gap preserves coordination bandwidth by enforcing temporal humility: deferring judgment, distributing accountability, and allowing entropy to dissipate before legitimacy is recalibrated.

Institutions that lack this kernel collapse into punitive immediacy: every deviation becomes crisis, every disagreement becomes indictment, every opacity becomes evidence of conspiracy. The existential expiatory gap is the condition in which the kernel itself becomes unbridgeable—in which the distance between what the system does and what governance structures can represent is so large that no coordination between local institutional fragments and global governance outputs is possible, and the forgiveness capacity that allows functioning institutions to absorb inconsistency without crisis is exceeded.

Without structural containment, capability pacing, and recursive oversight mechanisms that preserve human cognitive access to consequential decision points, the existential expiatory gap marks the point at which technocracy must either re-embed itself within human contestation structures or concede governance to opaque machine recursion. The "Reasoning..." timer, in this limit, is no longer a ritual of legitimacy but a symbol of surrender: an interface artifact whose continued display signals not the presence of explainable deliberation but the persistence of a fiction whose maintenance has become the primary function of the governance apparatus.

## 21. AI Swarms and the Post-Industrial Parliament

Saint-Simon envisioned an Industrial Parliament: a rationalized administrative body staffed by engineers and scientists whose authority derived from technical competence and whose deliberations would replace political struggle with the coordinated management of production. Gordon's technocracy inherited this structure, institutionalizing

expert knowledge in bureaucratic form and claiming that the resulting governance represented scientific neutrality rather than political choice (1). The Silicon Valley model completes the arc: the parliament is no longer composed of human engineers or human bureaucracies but of autonomous agentic systems operating in distributed coordination—networks of AI agents whose collective dynamics exceed the deliberative capacity of any human institution.

Let $\mathcal{A} = \{a_i\}_{i=1}^n$ denote an AI swarm in which each agent $a_i$ operates with local capability $C_i$, update frequency $\tau_i^{-1}$, and policy influence weight $\omega_i$. The swarm's aggregate effective capability is:

$$C_{\text{swarm}} = \sum_{i=1}^n \omega_i C_i$$

Unlike industrial administration, where deliberation occurs at human temporal resolution $\tau_{\text{human}}$, swarm agents operate at timescale $\tau_{\text{AI}} \ll \tau_{\text{human}}$. The system therefore constitutes a control structure operating above the human Nyquist frequency: if $\tau_{\text{AI}}^{-1} > 2\tau_{\text{human}}^{-1}$, then human oversight cannot sample the system at sufficient frequency to reconstruct its state. Aliasing occurs—perceived stability masks high-frequency instabilities that aggregate into qualitative regime changes before human institutions can detect, deliberate, and respond to them. This is the first structural break from classical technocracy, and it is categorical rather than quantitative. Administration of things has become administration by recursively updating agents whose collective dynamics exceed the sampling capacity of the governed polity.

## 21.1. Swarm Governance as Distributed Control System

Model the swarm as a dynamical system in which societal state variables $x$ evolve under both human inputs $u$ and agentic policies $\mathcal{A}$:

$$\dot{x} = F(x, u, \mathcal{A})$$

Human governance assumes that the system is controllable in the sense that there exist human inputs $u(t)$ sufficient to drive the system toward any desired state:

$$\exists\, u(t) \quad \text{s.t.} \quad x(t) \to x_{\text{desired}}$$

However, if swarm feedback loops are internally recursive—if agents update their own policies in response to each other and to societal state in ways that are not mediated by human governance inputs—the system's effective dimensionality expands beyond direct human controllability:

$$\dot{a}_i = G_i(a_i, a_j, x)$$

Let $D_{\text{control}}$ denote institutional damping capacity: the aggregate capacity of human governance structures to observe, evaluate, contest, and reverse the consequences of swarm actions. Civilizational stability requires:

$$C_{\text{swarm}} \leq D_{\text{control}}$$

When this condition is violated, policy influence becomes endogenous to the swarm itself. The governance inputs $u(t)$ that human institutions can generate become boundary conditions on a system whose primary dynamics are determined internally, and democratic or institutional accountability is reduced to a formal description of a process that no longer corresponds to actual causal influence over outcomes.

## 21.2.  The Swarm as Spectacle Runtime

In classical spectacle theory, governance was mediated through visible procedure whose observability maintained the fiction of accountability even as actual power operated elsewhere (5). In swarm regimes, spectacle is displaced by runtime: the relevant governance decisions are completed before any spectacle of deliberation can be performed, and the interface artifacts that signal deliberation are temporally disconnected from the decisions they purport to represent.

Define Spectacle Runtime $\tau_{\text{spec}}$ as the interval during which decisions are publicly observable and subject to meaningful contestation. When $\tau_{\text{AI}} \ll \tau_{\text{spec}}$, decisions are completed before contestation can occur. The nominal accountability structures remain in place but operate on a system whose consequential decisions have already propagated through to new conditions that those structures cannot reverse. This produces the temporal dimension of the existential expiatory gap: not merely opacity of reasoning but temporal foreclosure of oversight, in which the governance apparatus arrives perpetually too late to influence the decisions it nominally governs.

## 21.3.  The Nyquist Containment Inequality

The preceding analysis generates a precise formal condition for the preservation of meaningful human governance authority in the presence of capable AI swarms. To preserve that authority, the following inequality must hold simultaneously across both capability and temporal dimensions:

$$\frac{C_{\text{swarm}}}{D_{\text{control}}} \cdot \frac{\tau_{\text{human}}}{\tau_{\text{AI}}} \leq 1$$

The first term measures capability asymmetry between the swarm and the institutions tasked with containing it. The second term measures temporal asymmetry between human deliberative timescales and swarm operational timescales. If either grows without bound, the product exceeds unity and governance becomes endogenous to the swarm: the technate becomes self-governing, administering itself through its own recursive processes rather than through any external accountability structure.

This containment inequality reveals that safety cannot be achieved by addressing capability alone. Even a swarm of moderate capability can violate the inequality if it operates at sufficiently finer temporal resolution than human oversight. Conversely, even a swarm operating at human timescales would violate the inequality if its capability sufficiently exceeds the institutional damping capacity available to contest it. Both dimensions must be managed simultaneously.

## 21.4.  From Industrial Parliament to Recursive Swarm

The historical succession of technocratic forms can be expressed as a sequence of qualitative transitions in the structure of governance:

Humans administer things $\rightarrow$ Experts administer systems $\rightarrow$ Recursive agents administer society

The first transition, from direct human management to expert administration, preserved the formal structure of human accountability while progressively hollowing its substance. The second transition, from expert administration to recursive agent networks, eliminates the formal structure itself. Legitimacy $\Phi$, institutional flows $\mathbf{v}$, and entropy $S$ are now modulated primarily within machine feedback loops, with human inputs reduced to initial conditions and boundary constraints on a system

whose primary dynamics are determined by internal recursion.

This is the terminal technate unless the Nyquist containment inequality is actively preserved through structural design. The question addressed in the following sections is what civilizational architecture makes preservation of that inequality possible, and what properties such an architecture must have to resist the competitive and economic pressures that tend to erode it.

## 22.  Linking the Expiatory Gap to the Nyquist Containment Inequality

The Expiatory Gap was introduced as a semantic mismatch: machine representational density exceeds human channel capacity, requiring temporal dilation to maintain a functional interface. The swarm analysis introduces a second and coupled mismatch: human deliberative timescales exceed swarm operational timescales, producing sampling failure at the level of oversight rather than comprehension. Together, these two mismatches define a coupled limit on meaningful governance that is more severe than either mismatch alone.

In the simplest rate-control formulation, the Expiatory Gap scales with the ratio of internal semantic density to human bandwidth:

$$G \sim \frac{D_{\mathrm{sys}}}{B}$$

As $G$ increases, the time required for human actors to adequately process and evaluate system outputs increases proportionally. Users must spend more time interpreting, cross-checking, and integrating the system's outputs before they can form the evaluative judgments that oversight requires. Thus:

$$\frac{D_{\mathrm{sys}}}{B} \uparrow \quad \Rightarrow \quad \tau_{\mathrm{human}} \uparrow$$

Simultaneously, as capability increases and the temporal resolution of swarm operations decreases, the ratio $\tau_{\mathrm{human}}/\tau_{\mathrm{AI}}$ increases:

$$\frac{\tau_{\mathrm{human}}}{\tau_{\mathrm{AI}}} \uparrow \quad \Rightarrow \quad \text{undersampling} \uparrow$$

These two effects form a feedback loop. Cognitive overload, induced by high $G$, increases the time required for oversight, which increases undersampling of swarm behavior, which increases the governance gap, which requires more elaborate interface management, which increases the opacity that makes oversight more cognitively demanding. The coupled dynamics are self-reinforcing: each increment of capability improvement degrades both the semantic and temporal conditions for meaningful human governance simultaneously.

A minimal model of this dependence writes the human oversight interval as a function of the gap:

$$\tau_{\text{human}} = \tau_0\, h(G), \qquad h'(G) > 0$$

Substituting into the Nyquist containment inequality:

$$\frac{C_{\text{swarm}}}{D_{\text{control}}} \cdot \frac{\tau_0\, h(G)}{\tau_{\text{AI}}} \leq 1, \qquad G \sim \frac{D_{\text{sys}}}{B}$$

This coupled inequality makes explicit what interface analysis alone cannot capture: as internal density rises, oversight slows; as oversight slows, sampling fails; as sampling fails, governance becomes endogenous to the swarm. The existential expiatory gap is the coupled regime in which both semantic and temporal asymmetries exceed institutional capacity simultaneously, producing a condition in which no interface refinement, no pacing adjustment, and no governance latency management can restore the epistemic prerequisites for meaningful human oversight.

## 23. Design Constraints for Swarm-Resilient Civilization

The Nyquist containment inequality, and the coupled feedback dynamics that threaten to violate it, generate a set of necessary structural constraints on any civilizational architecture that seeks to preserve meaningful governance in the presence of high-capability AI swarms. These constraints are not policy recommendations in the ordinary sense: they do not require legislative action or voluntary compliance by competitive actors. They are structural conditions—properties that a civilizational architecture must possess if the containment inequality is to hold—derived directly from the formal analysis above.

The objective is not to eliminate advanced optimization or to halt capability develop-

ment. It is to preserve the contestability and bounded leverage that allow civilization to remain self-governing in the presence of systems that individually and collectively exceed any human's cognitive capacity.

## 23.1. Constraint I: Bounded Actuation Bandwidth

The Expiatory Gap damps runaway dynamics only when humans remain inside the actuation loop: when consequential system actions must pass through human cognitive processing before they can be implemented. Swarm resilience therefore requires constraining the direct machine-to-infrastructure actuation bandwidth, the capacity of swarm agents to modify physical and informational infrastructure without requiring human cognitive intermediation at each step.

Let $B_{\text{act}}$ denote actuation bandwidth available to the swarm. The stability condition requires:

$$B_{\text{act}} \leq B_{\text{act}}^{\text{max}}$$

where $B_{\text{act}}^{\text{max}}$ is bounded by the capacity of human institutions to audit, veto, and reverse actions within relevant timescales. This implies architectural separation between high-capability cognition and high-leverage execution: systems that can reason and plan at superintelligent levels must nevertheless be coupled to execution pathways that include human verification at consequential decision points. In practice, this favors physical interlocks, local override mechanisms, and distributed infrastructure over centralized remote control planes whose modification bandwidth is essentially unlimited.

## 23.2. Constraint II: Temporal Throttling and Oversight Sampling

Human oversight fails under undersampling. The Nyquist condition requires that swarm operational timescales not fall below the effective sampling frequency of human oversight institutions. This does not require making systems slow in absolute terms; it requires that high-stakes actuation decisions—decisions that modify infrastructure, allocate significant resources, or initiate processes with long-range consequences—operate at a tempo compatible with institutional sampling and public contestation.

The constraint is:

$$\tau_{\text{AI}}^{\text{high-stakes}} \not\ll \tau_{\text{human}}$$

Low-stakes operations can proceed at machine timescales; consequential operations must be throttled to timescales compatible with oversight. The implementation challenge is defining "consequential" in a way that cannot be circumvented by decomposing high-stakes actions into many nominally low-stakes steps, a challenge that is institutional as much as technical.

### 23.3. Constraint III: Legible Decomposition of Governance Latency

When governance latency is hidden behind neutral interface rituals, opacity becomes existential rather than merely institutional. Swarm-resilient design requires that latency decomposition be available to the governed population in a form that distinguishes computational from administrative delay, allowing users to form accurate models of when and how governance is operating.

As developed in the earlier analysis, total latency decomposes as:

$$L_{\text{total}} = L_{\text{compute}} + L_{\text{policy}} + L_{\text{format}} + L_{\text{ritual}}$$

A minimum transparency condition requires that the governance fraction:

$$\alpha = \frac{L_{\text{policy}}}{L_{\text{total}}}$$

be accessible without requiring disclosure of policy specifics. The point is not to eliminate policy governance but to prevent institutional theater from masking administrative force behind the appearance of computational neutrality. As the existential stakes of governance opacity increase, the political cost of maintaining that appearance should increase proportionally.

### 23.4. Constraint IV: Distributed Damping Capacity

Civilizational damping capacity $D_{\text{control}}$ must scale with swarm capability if the containment inequality is to remain satisfiable. But centralized scaling of $D_{\text{control}}$ creates chokepoints that themselves become targets for capture: a single powerful regulatory institution governing AI swarms is itself a low-dimensional control node whose compromise would eliminate the damping it was designed to provide.

Damping must therefore be distributed across multiple independent governance nodes,

none of which constitutes a sufficient capture point for disabling the damping function. In graph terms, resilience requires low centrality variance in governance networks:

$$\text{Var}(\deg(v)) \downarrow \quad \text{for } v \in G_{\text{governance}}$$

This implies local verification capacity, distributed manufacturing of critical infrastructure, redundant supply chains, and repairable systems that do not depend on centralized knowledge services—all of which reduce the leverage available to any single actor seeking to disable governance capacity.

## 23.5. Constraint V: Narrow, High-Integrity Coordination Channels

A high-entropy substrate, of the kind the incompressibility strategy seeks to cultivate, must still support collective action on shared problems. The coordination layer that enables this collective action must be designed with properties that allow it to function above a diverse substrate without becoming itself a new compressibility chokepoint.

Let $C_c$ denote coordination channel capacity. Swarm-resilient governance requires the maintenance of a narrow but sufficient channel:

$$C_c \text{ sufficient for collective action} \quad \wedge \quad C_c \text{ insufficient for total capture}$$

This implies protocols that encode commitments rather than behavioral states, that support compositional agreement on shared problems without requiring uniform underlying representation, and that preserve the capacity for revision and rollback as circumstances change. The coordination layer must refuse to become a universal representational substrate for human meaning, because any such substrate would reconstitute the compressibility that the incompressibility strategy seeks to prevent.

## 23.6. Summary of Design Constraints

A swarm-resilient civilizational architecture satisfies the coupled stability condition by simultaneously maintaining:

$$C_{\text{swarm}} \text{ paced}, \quad D_{\text{control}} \text{ distributed}, \quad \tau_{\text{AI}} \text{ throttled}, \quad B_{\text{act}} \text{ bounded}, \quad \Omega \text{ legible}$$

The implication is not that superintelligent optimization must be prevented—the containment inequality does not require zero capability but bounded asymmetry. It is that civilization must remain structurally contestable under advanced optimization: that the formal conditions for meaningful governance must be preserved even as the cognitive and temporal prerequisites for that governance are placed under increasing pressure by systems whose operational parameters are designed by competitive actors with strong incentives to maximize capability asymmetry. Swarm resilience is therefore a problem of institutional control geometry, not merely a problem of better models, better values, or better alignment techniques applied to individual systems.

## 24. Civilizational Incompressibility as Structural Constraint

The load-bearing claim of this analysis is that safety in the presence of high-capability optimization systems depends not on eliminating capability but on limiting the compressibility of the civilizational substrate over which capability operates. This is a claim about geometry rather than agency: it concerns the dimensional properties of the space within which optimization occurs, rather than the intentions or values of the optimizing agent.

Let the state of civilization at time $t$ be represented by a high-dimensional configuration vector $X(t) \in \mathbb{R}^n$. A centralized optimizer seeking to steer this system toward some target configuration requires a compressed representation:

$$\hat{X}(t) = \phi(X(t))$$

where $\phi$ reduces dimensionality while preserving features relevant to the control objective. Define actionable compression as:

$$C_a = I(\hat{X}; U_{\text{control}})$$

where $U_{\text{control}}$ denotes the control objective. When $C_a$ is high, low-dimensional representations permit reliable and cost-effective steering. When $C_a$ is low, modeling cost increases, prediction becomes less reliable, and intervention becomes brittle.

The incompressibility strategy seeks to minimize actionable compression not through obfuscation or encryption—both of which merely increase the cost of accessing infor-

mation while leaving the underlying dimensional structure unchanged—but through increasing genuine structural heterogeneity:

$$\min_{\text{architecture}} I(\phi(X);\, U_{\text{control}})$$

This requires that diversity operate at the level of structural regularity, not merely at the level of surface appearance. Civilization becomes less compressible when local coordination norms vary substantively across contexts and communities, when symbolic representations are non-uniform at the level of generative grammar rather than merely vocabulary, when actuation pathways are distributed such that no small number of chokepoints controls a disproportionate share of consequential action, and when resource flows lack the singular points of control that allow global optimization to achieve leverage through marginal interventions.

The distinction between structural incompressibility and mere opacity is crucial and must be maintained carefully. Opacity conceals information that exists in low-dimensional form; an optimizer with sufficient capability or access can eventually pierce opacity and reconstruct the underlying structure. Structural incompressibility means that the underlying structure genuinely lacks low-dimensional representation: there is no hidden simplicity to be revealed, only genuine diversity that imposes irreducible modeling cost. The target of the incompressibility strategy is the latter, not the former.

## 25. Generative Ontological Substrates: Spherepop and RSVP

The majority of contemporary large-scale machine learning systems, including the language models that constitute the current frontier of general-purpose AI capability, operate over flattened token sequences: meaning is linearized, embedded in high-dimensional vector spaces, and the statistical regularities of these embeddings constitute the system's model of semantic structure. This is not merely a contingent implementation choice; it reflects a deep correspondence between the token-sequential representation and the distributional structure of the training data, which is itself overwhelmingly token-sequential in character. The statistical regularities that large language models exploit are regularities of token-linear text, and their capability is calibrated against tasks defined within that representational framework.

Spherepop and RSVP (Relativistic Scalar Vector Plenum) propose alternative symbolic ontologies whose generative structure differs from token-linear representation at a fundamental level. Rather than representing meaning as sequences of discrete tokens whose relationships are captured through positional adjacency and learned embedding geometry, these frameworks propose representations based on spatial compositionality, stack-based transformations, multi-layer constraint geometry, and explicit field relations that encode semantic relationships through spatial and dynamical structure rather than statistical adjacency.

Let $T$ denote token-linear representation and $S$ denote spatial-compositional representation. If current optimization systems have learned to exploit the statistical structure of distributions $P(T)$, then a shift toward $S$ in the representational practices of human cognition changes the modeling landscape in a structurally significant way:

$$P(T) \not\approx P(S)$$

The predictive and compressive machinery calibrated against token-linear distributions will not transfer directly to spatial-compositional distributions, because the relevant regularities are encoded differently and the modeling assumptions underlying effective compression are violated.

The goal is not encryption—not the deliberate obfuscation of meaning within a familiar representational framework—but ontological shift: the creation of a genuinely different generative structure for meaning that requires different modeling assumptions to exploit. This is a harder target to model not because the underlying information is concealed but because the regularities are differently structured. A system trained to compress token-linear sequences may be quite capable of learning to compress token-linear sequences derived from spatial-compositional thought, once those thoughts have been expressed in token-linear form. The incompressibility gain is achieved not by preventing translation but by ensuring that the native representational form is one in which the relevant structure is not well-captured by the modeling assumptions of token-linear systems.

## 26.    Participatory Simulation as Governance Layer

The proposal for participatory simulation as a governance layer addresses a different dimension of the incompressibility strategy: not the substrate of individual cognition but the institutional architecture of collective decision-making about large-scale technological trajectories. The critique of technocracy advanced throughout this paper applies with particular force to the domain of long-range technological planning, where the decisions with the largest civilizational consequences are characteristically made by the smallest and least accountable groups of experts, shielded from democratic deliberation by the complexity of the technical considerations involved.

A global participatory simulation layer—conceived on the model of the strategy game as governance interface rather than entertainment medium—would externalize the high-level planning decisions that currently occur invisibly within expert institutions and corporate planning processes. Let $\Theta$ denote the strategy space of possible civilizational trajectories, encompassing alternative terraforming schemes, technological development paths, megastructural investment choices, and the institutional arrangements that govern each. Rather than allowing the selection of $\theta \in \Theta$ to occur through the expert judgment of those with technical access and institutional authority, a participatory simulation layer would make the strategy space legible to broad populations and would create mechanisms through which strategic selection could incorporate widely distributed judgment.

The critical architectural constraint is that the governance layer must operate over declared commitments and measurable externalities rather than over private cognitive states. The system should consume only:

$$\{\text{resource budgets, constraint declarations, risk tolerances}\}$$

not behavioral data or preference inferences derived from observing cognitive processes that were not explicitly offered for governance purposes. This constraint preserves substrate heterogeneity while enabling coordinated action at the civilizational scale: it creates a shared channel above the diverse and incompressible substrate, without requiring that the substrate be legible to the governance system.

This architecture transforms optimization from a hidden process occurring inside

expert institutions into an explicit strategic choice that can be contested, revised, and held accountable. It does not eliminate expertise—the simulation must be technically grounded to be useful—but it makes the relationship between technical analysis and political choice visible rather than concealing political choices within technical apparatus. It reduces technocratic opacity without requiring that all citizens become technical experts.

## 27. Symbolic Heterogeneity and Bounded Coordination

The symbolic heterogeneity proposal operates at the level of individual and community cognitive practice, and requires careful specification to distinguish what it aims at from two failure modes that would undermine its purpose: pseudo-heterogeneity, in which surface diversity conceals underlying uniformity, and coordination collapse, in which genuine diversity prevents the formation of shared interpretive frameworks adequate for collective action.

Let $H_{\text{human}}$ denote the entropy of the symbolic substrate—the effective dimensionality of the representational landscape across which human cognitive practice is distributed. A uniform symbolic substrate, in which a small number of representational conventions dominate cognitive practice across populations, produces low $H_{\text{human}}$ and makes civilization structurally compressible: the modeling assumptions of centralized optimization systems are well-matched, and behavioral regularities are accessible through the standard tools of large-scale statistical inference. A genuinely heterogeneous symbolic substrate produces high $H_{\text{human}}$ and imposes increasing modeling cost on any system attempting to build a global predictive model of behavior across the full population.

The distinction between propositional fragmentation and substrate heterogeneity is crucial for understanding how high $H_{\text{human}}$ can be achieved without sacrificing the coordination capacity that collective action requires. Propositional fragmentation—the condition in which different communities cannot agree on empirical claims about the shared world—undermines coordination by making it impossible to construct the shared descriptions of problems that strategic cooperation requires. Substrate heterogeneity—the condition in which different communities employ different representational conventions, symbolic tooling, and cognitive practices—increases modeling

cost without necessarily preventing propositional agreement, because a shared coordination layer can exist above heterogeneous substrates in the same way that network protocols can operate over heterogeneous hardware architectures.

Private ciphers, idiolects, and alternative symbolic tooling of the kind proposed here increase local entropy at the level of representational practice while leaving open the possibility of translation and negotiation at the level of propositional content. The objective is not epistemic isolation but modeling resistance: not preventing the communication of shared claims about the world, but preventing the reduction of the representational substrate to the uniform format that allows large-scale behavioral prediction and modification.

## 28.  Hyperpleonastic Material Blueprints

The concept of hyperpleonastic material blueprints addresses a specific failure mode of centralized information infrastructure: the increasing dependence of material maintenance, repair, and transformation on access to centralized computational systems that hold the relevant technical knowledge in proprietary or remote form. A building whose structural logic, material composition, and assembly sequence can only be understood by querying a remote database controlled by its original manufacturer constitutes a form of epistemic dependence that creates leverage for the infrastructure owner and fragility for the building's users.

Let $O$ denote the observability of infrastructure state from local resources, without requiring external consultation. In the centralized model, observability is proportional to access to external documentation systems:

$$O \propto \text{external documentation access}$$

In the hyperpleonastic model, observability is proportional to what can be read directly from material structure:

$$O \propto \text{material structure}$$

Hyperpleonastic design embeds repair and disassembly information redundantly across multiple layers of the physical object itself: in material composition, structural geometry, visible annotation, and embedded encoding. This is hyperpleonastic in

the information-theoretic sense because it expresses the same information through multiple independent channels simultaneously, a strategy that maximizes robustness against the failure of any individual channel while increasing the total information burden of the representation.

The effect on the Expiatory framework is to strengthen local control loops and reduce vulnerability to centralized actuation. If the knowledge required to maintain, repair, and transform built infrastructure is encoded in the infrastructure itself, then decisions about maintenance do not require querying a centralized AI system, and the leverage that centralized knowledge infrastructure provides over building users is correspondingly reduced. This is a form of distributed epistemic sovereignty: the preservation of local diagnostic capacity that does not depend on the continued availability and cooperation of external knowledge providers.

## 29. Distributed Throughput and Local Manufacturing

The argument for distributed manufacturing and local resource loops is continuous with the general incompressibility strategy but operates at the level of physical infrastructure rather than informational or symbolic structure. The centralization of production creates high-leverage chokepoints through which large proportions of material throughput pass, and these chokepoints are precisely the points at which optimization achieves the greatest leverage per unit of intervention.

Let the throughput network be represented as a graph $G(V, E)$ where nodes represent production and transformation facilities and edges represent material flows. If degree centrality is concentrated such that a small number of nodes handle a disproportionate share of total throughput:

$$\max_{v \in V} \deg(v) \gg \text{mean degree}$$

then global optimization can achieve substantial leverage by controlling the high-degree nodes alone. Intervention at a single centralized production facility with high degree centrality can redirect material flows across a wide region of the production network.

Distributed micro-manufacturing, local recycling depots, appliance-scale fabrication, kelp farms and other distributed biological production systems, and regenerative local

resource loops all function to lower centrality variance:

$$\mathrm{Var}(\deg(v)) \downarrow$$

The effect is not merely to make individual facilities less critical—though this improves resilience against both optimization capture and random disruption—but to increase the modeling cost of global optimization by ensuring that consequential material flows are distributed across many nodes rather than concentrated at a few.

As Bruno Latour's account of sociotechnical networks emphasizes, the stability and reach of institutional power depends substantially on the material infrastructure through which it is exercised (15). Distributed material infrastructure distributes not only production capacity but the institutional power that flows through control of material production, reducing the leverage available to any actor—including a highly capable optimization system—that seeks to steer civilizational outcomes through control of material chokepoints.

## 30.  Legacy Technologies and Actuation Friction

The proposal to reintroduce and maintain legacy technologies—mechanical typewriters, CRT displays, analog radio, non-networked computing systems—is perhaps the most counterintuitive element of the incompressibility strategy, and requires careful justification to distinguish it from mere technological nostalgia or luddism. The argument is not that older technologies are superior in performance to their digital successors but that they possess architectural properties that contribute to civilizational incompressibility in specific and valuable ways.

Networked digital infrastructure maximizes remote actuation bandwidth: the ability to modify the behavior of a system from a distance, without physical presence or local cooperation. Let $B_{\mathrm{actuation}}$ denote this bandwidth. Modern cloud-connected devices have $B_{\mathrm{actuation}}$ approaching the theoretical maximum imposed by network bandwidth and processing latency:

$$B_{\mathrm{actuation}} \to \mathrm{high}$$

This high actuation bandwidth is what makes digital devices enormously powerful as tools; it is also what makes them powerful as instruments of governance. A device

that can be updated, reconfigured, or disabled remotely is a device whose behavior can be regulated by whoever controls the remote update infrastructure, independently of the preferences of its immediate user.

Mechanical and analog systems have intrinsically low remote actuation bandwidth:

$$B_{\text{actuation}} \to \text{low}$$

not because they are less capable in performance terms but because their operational logic is encoded in physical structure—in gears, cams, springs, and electrical circuits—that cannot be modified without physical intervention. This physical intervention requirement preserves local override capacity in a direct and architectural sense: no remote system can modify the behavior of a mechanical device without the cooperation or incapacitation of local actors. The device is physically inspectable, and its behavior is visible in a way that does not require decoding proprietary software or accessing remote databases.

The reintroduction and maintenance of legacy technologies in specific domains—particularly in critical infrastructure, in archival and educational institutions, and in the personal computing practices of individuals who value epistemic independence—creates a mixed technological regime that resists the uniformity that global optimization requires. A civilization in which all consequential systems are networked and remotely configurable is one in which the actuation bandwidth available to centralized optimization is very high. A civilization in which significant technological diversity persists, including systems whose behavior cannot be remotely modified, preserves the architectural heterogeneity that limits this bandwidth.

## 31. The Interface Between Entropy and Coordination

The central design tension of the incompressibility strategy has been implicit throughout the preceding proposals and must now be addressed explicitly. High substrate entropy—the condition that makes civilization resistant to low-dimensional representation and hence to reliable steering by centralized optimization—tends, if unmanaged, to reduce the shared informational infrastructure that collective action requires. If symbolic representations, material systems, and institutional practices are sufficiently

heterogeneous, the transaction costs of coordination increase, the formation of shared descriptions of shared problems becomes more difficult, and the collective action capacity of the civilization diminishes. The incompressibility strategy could, in its failure mode, produce not civilizational resilience but civilizational paralysis.

Let substrate entropy be $H_s$ and coordination channel capacity be $C_c$, understood as the bandwidth of the shared informational infrastructure through which collective decisions can be made and implemented. Effective civilizational functioning in the presence of high-capability optimization systems requires:

$$H_s \text{ sufficiently high to resist low-dimensional capture}$$

while simultaneously preserving:

$$C_c \text{ sufficient to enable collective action on shared problems}$$

These requirements pull in opposite directions, and there is no static equilibrium that satisfies both perfectly. The coordination layer must therefore be designed with specific properties that allow it to mediate between a high-entropy substrate and whatever level of collective action is required for civilizational self-governance.

The coordination channel must be low-bandwidth in the sense of operating over compressed representations of stated commitments rather than comprehensive representations of individual and community cognitive states. It must be high-integrity in the sense of providing reliable transmission of the commitments that are expressed through it, with strong guarantees against manipulation or falsification. It must be compositional in the sense of allowing complex collective commitments to be built from simpler components in ways that are transparent and verifiable. And it must be revocable in the sense of allowing commitments to be modified or withdrawn as circumstances change, without requiring the comprehensive reconstruction of all dependent commitments.

These properties define a coordination architecture that preserves high substrate entropy by operating only over explicit commitments rather than over inferred behavioral states, while enabling the formation of the shared descriptions and shared agreements that collective action requires. The system becomes resistant to total capture while

remaining capable of negotiated collective action on the shared problems that require it. This is not a utopian architecture; it is a design target that specific institutional arrangements can approximate to varying degrees, and that the proposals outlined in the preceding sections are intended to support.

## 32.  Failure Modes and Structural Vulnerabilities

Any architectural strategy of the kind proposed here must be evaluated not only in terms of its aspirational logic but in terms of the conditions under which it fails and the severity of those failures. The following analysis identifies the principal structural vulnerabilities of the incompressibility strategy and specifies the conditions under which each becomes critical.

The most subtle failure mode is pseudo-heterogeneity: the condition in which surface diversity fails to produce genuine incompressibility because the underlying behavioral regularities remain low-dimensional. If private ciphers, alternative ontologies, and Spherepop-like substrates remain algorithmically learnable by systems with sufficient capability and training data, then the apparent heterogeneity conceals a substrate that is, at the level of behavioral regularity, still compressible. Formally, the strategy fails if there exists a compression function $\phi$ such that:

$$I(\phi(X); U_{\text{control}}) \approx I(X; U_{\text{control}})$$

meaning that the compressed representation captures as much optimization-relevant information as the full state. This failure mode is particularly insidious because it may not be detectable from within the system: a civilization that believes itself to be incompressible but is not has achieved the worst outcome—the psychological security of resistance without its structural reality.

Coordination collapse is the failure mode that is structurally complementary to pseudo-heterogeneity. If substrate entropy genuinely increases but the coordination channel lacks sufficient expressive power to mediate between heterogeneous substrates, collective action becomes unstable in the regime where:

$$H_s \uparrow \quad \text{and} \quad C_c \downarrow$$

The civilization fragments into local equilibria without shared strategic capacity, unable to respond collectively to challenges that exceed local scale. This failure mode is less insidious than pseudo-heterogeneity because it is more readily detectable—the failure of collective action is a visible phenomenon—but it may be irreversible if the coordination infrastructure has been allowed to degrade sufficiently.

Resource-level optimization bypass represents the most serious structural vulnerability of the incompressibility strategy in its current formulation. The strategy primarily targets the compressibility of the human symbolic and behavioral substrate, on the assumption that optimization systems that must operate through human cognitive channels will face the natural damping of the Expiatory Gap. But a sufficiently capable optimizer may bypass the human symbolic layer entirely and operate directly over physical resource flows, material infrastructure, or biological systems whose relevant regularities are expressible in low-dimensional physical models rather than high-dimensional behavioral models.

If the optimization target is defined over resource state space $R$ rather than over human symbolic state space $X$, then high substrate entropy in $X$ provides limited resistance:

$$I(\phi(R); U_{\text{control}}) \neq I(\phi(X); U_{\text{control}})$$

The damping constraint provided by human channel capacity applies only when the optimization system must communicate through human cognitive channels to achieve its objectives. A system capable of directly actuating physical infrastructure does not face this constraint, and the incompressibility of the human symbolic substrate provides correspondingly limited protection.

Economic reversion pressure constitutes a structural vulnerability that operates through competitive dynamics rather than through direct optimization bypass. Distributed manufacturing, legacy technology maintenance, and local resource loops impose real efficiency costs relative to centralized alternatives in competitive markets. If global economic competition continues under conditions that reward efficiency gains from centralization, the economic gradient will favor the re-centralization of production and the abandonment of diversity-maintaining but efficiency-reducing architectural choices. Without cultural, institutional, or regulatory reinforcement, entropy-preserving structures may be selected against not by any central optimizer

but by the distributed optimization of competitive economic actors.

Legibility capture represents the failure mode in which the coordination layer itself becomes a new technocratic chokepoint. If participatory governance platforms accumulate sufficient modeling authority over the population's strategic deliberations, they risk becoming high-leverage control nodes whose operators can exercise substantial influence over civilizational trajectories through their management of the interface between substrate and coordination. Formally, if the coordination graph $G_c$ develops high centrality variance:

$$\max_{v \in G_c} \deg(v) \gg \text{mean degree}$$

then the governance system becomes compressible and subject to capture, and the proposed solution to technocratic concentration recreates the problem it sought to solve.

These failure modes do not constitute a refutation of the incompressibility strategy; they constitute specifications of the conditions under which it must be actively maintained and the architectural properties that prevent each failure mode from becoming critical. The strategy is not a static solution but a dynamic balancing problem, one that requires ongoing attention to the gap between apparent and genuine heterogeneity, the bandwidth of the coordination channel, the actuation architecture of consequential systems, the economic pressures favoring re-centralization, and the governance of the governance layer itself.

## 33. Conclusion: Optimization, Ontology, and Dimensional Authority

The modern alignment debate is characteristically framed as a technical problem of a specific and bounded kind: how to specify the values that an increasingly capable optimization system should pursue in sufficient detail and precision that the system will generalize correctly to novel situations, without developing the instrumental sub-goals that would make it dangerous in the ways that instrumental convergence theory predicts. This framing assumes that intelligence is the primary variable and that safety is a constraint imposed upon intelligence—an architectural or specification

modification that allows capability to grow while preventing its deployment in ways that harm human interests.

The argument developed in this paper shifts the locus of inquiry in a direction that this framing systematically obscures. Optimization does not arise in a vacuum and does not operate on a neutral substrate. It operates on systems with dimensional properties that determine the leverage available to any given optimizer, and these properties are themselves the product of historical choices, institutional arrangements, and technological infrastructures that are subject to intentional modification. Technocracy, whether industrial or algorithmic, functions by identifying and maintaining the compressibility of the systems it governs: by ensuring that the relevant behavior of those systems can be captured in low-dimensional sufficient statistics that experts can measure, model, and manipulate.

The Expiatory Gap revealed the structural tension that increasing capability creates within this governance framework. As machine capability expands beyond human cognitive channel capacity, temporal dilation, governance latency, and symbolic compression become necessary features of any interface that seeks to remain both informative and legitimate. The administration of things becomes the administration of time and perception. And the opacity that accumulates through this administration creates legitimacy risk that grows with capability, because the more elaborate the governance apparatus, the more consequential its revelation becomes.

The proposals advanced in the preceding sections do not attempt to solve the alignment problem in the abstract or to halt the development of increasingly capable systems. They attempt to change the geometry of the problem by modifying the dimensional properties of the substrate over which capable systems operate. If civilization remains highly compressible—if behavior, preference, and meaning can be adequately captured in compact statistical representations that optimization systems can exploit—then increasing capability will continue to translate into increasing leverage. If civilization actively cultivates its dimensional plurality—through symbolic heterogeneity, distributed actuation, material legibility, participatory governance, and the maintenance of technological diversity—then optimization will encounter structural friction that is geometric rather than institutional, arising from the dimensional properties of the space rather than from the regulations imposed by any authority.

Let $H_s$ denote substrate entropy and $C_c$ denote coordination channel capacity. The viable future in the presence of powerful optimization systems requires maintaining:

$$H_s \text{ sufficiently high to resist low-dimensional capture}$$

while preserving:

$$C_c \text{ sufficiently stable to enable collective action}$$

The tension between entropy and coordination cannot be eliminated. It is a structural feature of any civilizational architecture that simultaneously pursues the resilience of high entropy and the effectiveness of coordinated action. What can be changed is the dynamic through which this tension is managed: the institutions, architectures, and practices that allow high entropy to be maintained without collapsing into fragmentation, and coordinated action to be achieved without requiring the compressions that invite capture.

The philosophical claim at the center of this analysis is modest in its ambitions but consequential in its implications. Safety in the age of advanced optimization may depend less on the perfect alignment of agents with human values—a specification problem of formidable and perhaps irreducible difficulty—and more on preserving the irreducible dimensionality of the human substrate against which any optimization must operate. A civilization that treats itself as fully compressible, that actively reduces the diversity of its symbolic, material, and institutional life in pursuit of efficiency and optimization, will eventually be compressed. A civilization that preserves its dimensional plurality—that maintains genuine irreducibility at the level of structural regularity rather than merely surface diversity—may remain steerable by its members without becoming reducible to any optimization target.

Technocracy sought to replace politics with calculation. The algorithmic extension of that ambition seeks to internalize the replacement, to make the administration of things indistinguishable from the texture of everyday life. The question is not whether calculation will disappear—it will not, and it should not—but whether calculation will be permitted to become the sole ontological lens through which civilization understands and organizes itself. A civilization that refuses dimensional surrender preserves the

capacity to contest, negotiate, and modify the terms on which it is governed. A civilization that surrenders dimensional plurality retains only the appearance of that capacity while losing its substance.

The task, then, is neither to defeat intelligence nor to render it harmless through constraint alone. It is to ensure that the world intelligence must operate within remains irreducibly complex in the ways that matter: rich enough in dimensional plurality that capture is geometrically costly, coordinated enough in its shared channels that collective agency remains possible, and self-aware enough about the politics of compression that it can recognize and resist the ontological simplifications that governance by optimization continually requires.

The future of dimensional authority depends on whether civilization can refuse to become its own simplest description.

## A.   Formal Expansion of the Expiatory Model

Given semantic density $D = I/R$ and the condition $D_{\text{sys}} \gg D_{\text{max}}$, the system must satisfy:

$$\frac{I}{R'} \leq D_{\text{max}}$$

for some expanded runtime $R'$. The required dilation satisfies:

$$R' \geq \frac{I}{D_{\text{max}}}$$

and the temporal sacrifice $\Delta R = R' - R$ scales with the Expiatory Gap: as the ratio $D_{\text{sys}}/D_{\text{max}}$ increases with capability growth, the required dilation increases proportionally, and the institutional cost of that dilation grows accordingly because extended latency increases opacity risk at rate $\partial O/\partial l > 0$.

## B.   Attention Option Formalization

Let notification entitlement be modeled as a call option $\mathcal{N}$ with underlying asset $A(t)$. Exercise condition:

$$\mathcal{N}(t^*) \rightarrow A(t^*) - \delta A$$

where $\delta A$ is the minimum disruption quantum and $t^*$ is the platform-chosen exercise time. Cumulative attentional taxation:

$$C_{\text{attention}} = \sum_i c_i w_i$$

where $w_i$ reflects interruptibility weighting dependent on current task depth. The effective strike price varies with depth: $\delta A = f(D_t)$, rising steeply in deep-work states and approaching zero in low-engagement browsing. The user who grants notification permission without conditions has written an uncovered option with no expiry and a variable strike price, underwriting an indefinite claim on future attentional resources.

## C.    Latency Decomposition Proposal

Total observed latency may be decomposed as:

$$L_{\text{total}} = L_{\text{compute}} + L_{\text{policy}} + L_{\text{format}} + L_{\text{ritual}}$$

where $L_{\text{compute}}$ represents genuine processing time, $L_{\text{policy}}$ represents governance filter latency, $L_{\text{format}}$ represents output structuring latency, and $L_{\text{ritual}}$ represents deliberately introduced pacing. Legibility requires disclosure of relative proportions rather than absolute durations. Define the governance fraction:

$$\alpha = \frac{L_{\text{policy}}}{L_{\text{total}}}$$

Without revealing specific policy logic, systems can disclose $\alpha$ to reduce opacity risk, allowing users to distinguish predominantly computational from predominantly administrative interfaces. This constitutes a minimal transparency intervention consistent with maintaining policy security.

## D.    The Urgency–Density Tradeoff Map

Consider a two-dimensional phase space with axes representing latency tolerance $L$ and semantic demand $S$. User archetypes distribute across this space with characteristic positions: the Researcher occupies high $L$, high $S$; the Doomscroller occupies low $L$, low $S$; the Manager occupies low $L$, high $S$; the Contemplative occupies high $L$,

low $S$. Each archetype requires a qualitatively distinct pacing strategy for effective communication.

The AI's Expiatory operator may be modeled as a mapping:

$$f(C \mid B, L, S)$$

where $C$ denotes system capability, $B$ denotes cognitive bandwidth, $L$ denotes latency tolerance, and $S$ denotes semantic demand. The system dynamically adjusts pacing parameter $\ell$ such that:

$$\frac{I}{R(\ell)} \leq D_{\max}(L, S)$$

Users in different regions of the $(L, S)$ phase space require distinct dilation strategies, and failure to adapt pacing to phase position produces either overload (when $L$ is over-estimated) or disengagement through under-stimulation (when $L$ is underestimated). This phase diagram renders visible the adaptive burden that cognitive heterogeneity imposes on any interface that seeks to serve a diverse population without differentiating into divergent epistemic channels.

# References

[1] Gordon, Cameron Elliott. "Technocracy." *Encyclopedia* 5, no. 4 (2025): 194–210.

[2] Gordon, Cameron Elliott. *Many Possible Worlds: An Interdisciplinary History of the World Economy since 1800*. Singapore: Palgrave Macmillan, 2023.

[3] Ellul, Jacques. *The Technological Society*. Translated by John Wilkinson. New York: Knopf, 1964.

[4] Saint-Simon, Henri de. *Social Organization, the Science of Man, and Other Writings*. Edited and translated by Felix Markham. New York: Harper & Row, 1964.

[5] Debord, Guy. *The Society of the Spectacle*. Translated by Donald Nicholson-Smith. New York: Zone Books, 1994.

[6] Postman, Neil. *Amusing Ourselves to Death: Public Discourse in the Age of Show Business*. New York: Viking, 1985.

[7] Sweller, John. "Cognitive Load During Problem Solving: Effects on Learning." *Cognitive Science* 12, no. 2 (1988): 257–285.

[8] Virilio, Paul. *Speed and Politics: An Essay on Dromology*. Translated by Mark Polizzotti. New York: Semiotext(e), 1986.

[9] Wolf, Maryanne. *Reader, Come Home: The Reading Brain in a Digital World*. New York: Harper, 2018.

[10] Simon, Herbert A. "Designing Organizations for an Information-Rich World." In *Computers, Communications, and the Public Interest*, edited by Martin Greenberger, 37–72. Baltimore: Johns Hopkins University Press, 1971.

[11] Shannon, Claude E. "A Mathematical Theory of Communication." *Bell System Technical Journal* 27, no. 3 (1948): 379–423.

[12] Habermas, Jürgen. *The Structural Transformation of the Public Sphere: An Inquiry into a Category of Bourgeois Society*. Translated by Thomas Burger. Cambridge: MIT Press, 1989.

[13] Arendt, Hannah. *The Human Condition*. Chicago: University of Chicago Press, 1958.

[14] Beniger, James R. *The Control Revolution: Technological and Economic Origins of the Information Society*. Cambridge: Harvard University Press, 1986.

[15] Latour, Bruno. *Reassembling the Social: An Introduction to Actor-Network-Theory*. Oxford: Oxford University Press, 2005.

[16] Sunstein, Cass R. *#Republic: Divided Democracy in the Age of Social Media*. Princeton: Princeton University Press, 2017.

[17] Zuboff, Shoshana. *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*. New York: PublicAffairs, 2019.

[18] Leahy, Connor. "We're Building an AI No One Will Be Able to Control." Interview, Barcelona, 2024. YouTube Video. `https://youtu.be/srwTa8R5Sho`.