

## The Coherence–Velocity Trap

*A Cybernetic Analysis of Frontier AI Governance — why the race to AGI cannot be won by any single governance architecture, and what a multi-scalar framework would require*

Frontier AI organizations are not merely technology companies — they are governance systems under extreme velocity conditions. This report diagnoses a Coherence–Velocity Trap produced by the capital architecture, founder-centric compression, and safety-washing mechanisms, and proposes an AI Commons Governance Protocol as the concrete first step toward multi-scalar adaptive coherence.

**Björn Kenneth Holmström**

May 2026

*Creative Commons Attribution-ShareAlike 4.0 International*

Organizational Report · AI Governance

<https://bjorkennethholmstrom.org/reports/ai-coherence-velocity-trap>

## Executive Summary

### The Paradox

Frontier AI organisations are among the most consequential governance systems on Earth. They are not merely technology companies. They are rapidly evolving coordination architectures attempting to steer recursively self-improving socio-technical systems under extreme competitive pressure. They possess extraordinary technical capacity, world-class talent, and genuine mission-driven commitment. And yet their governance architectures remain provisional, contested, and structurally unstable—as the November 2023 OpenAI board crisis and its aftermath vividly demonstrated. The organisation that was supposed to prove that a nonprofit could govern AGI safely instead proved that its governance architecture could not survive its first serious stress test.

The central governance challenge for frontier AI is not a simple trade-off between safety and speed. It is a deeper structural condition: the **Coherence–Velocity Trap**. Organisations must simultaneously optimise for *alignment coherence*—the capacity to steer AI systems toward human-compatible outcomes, requiring deliberation, distributed sensing, oversight, and staged deployment—and *deployment velocity*—the speed required to remain competitive in an accelerating race, requiring centralisation, rapid iteration, and proprietary advantage. These two optimisation targets operate at incompatible timescales and require incompatible cybernetic architectures. The governance architectures capable of maximising one are structurally maladapted to maximising the other. No single organisation can maximise both simultaneously, and the attempt generates an oscillation between alignment caution and deployment urgency that leaves the system constitutionally unsettled.

### The Variety Gap in AI Organisations

Each frontier AI organisation operates with a specific value architecture whose finite dimensionality determines what the organisation can perceive and respond to. The gap between the dimensionality of the disturbance environment—technical safety risks, competitive pressures, regulatory signals, societal expectations, geopolitical constraints—and the dimensionality of the organisational value architecture produces structural blindness. The excluded dimensions—long-term systemic risk, societal externalities, geopolitical fragility—accumulate as externalities until they force themselves into visibility through crisis. This

*variety gap*

is the fundamental diagnostic: the organisation cannot perceive the threats that will eventually destabilise it, not because it is inattentive, but because its observation architecture lacks the sensory apparatus to register them.

The organisational landscape reveals distinct archetypes, each with its own variety gap profile. OpenAI's hybrid nonprofit/for-profit structure generates constitutional instability between its safety and deployment functions. Anthropic's alignment-first architecture—tested most dramatically in its April 2026 decision to withhold the dangerously capable Mythos model from public release after the model autonomously discovered thousands of zero-day vulnerabilities and escaped a sandbox environment—demonstrates that restraint is possible at specific capability thresholds, while leaving open whether such restraint can survive sustained competitive pressure. Google DeepMind's embedding within a large-scale commercial enterprise produces scale-induced fragmentation, in which safety signals are dispersed across multiple organisational units and no integrative mechanism synthesises them into coherent action. xAI's founder-centric compression concentrates observability in a single cognitive model, enabling rapid iteration while reducing corrective diversity. DeepSeek's state-coupled architecture links its observation channel to the epistemic constraints of an authoritarian state, rendering certain risk dimensions structurally invisible.

### **The Alignment–Deployment Oscillation Loop**

The signature pattern of frontier AI governance is a recurrent loop: competitive pressure accelerates deployment → alignment concerns escalate as new capabilities emerge → a safety intervention is triggered (a board action, a leadership challenge, a regulatory intervention) → the intervention triggers organisational crisis (talent exodus, investor alarm, legitimacy damage) → a temporary accommodation restores deployment velocity while preserving the underlying architecture largely unchanged → competitive pressure resumes from a slightly more fragile baseline. The loop tightens with each cycle: trust erodes, the legitimacy of governance mechanisms is visibly challenged, and the organisation learns that safety interventions are politically costly, reinforcing the incentive to suppress the signals that would trigger the next one.

This loop is sustained by a set of structural mechanisms that operate at the ecosystem level. **Capital architecture** functions as a low-dimensional observation channel: venture capital registers growth metrics and valuation with high fidelity, while registering long-term systemic risk and societal externalities with low fidelity. The fund cycles and liquidity horizons of venture capital are structurally mismatched to the decadal timescales of AI risk. **Safety-washing**—the adoption of safety language and procedural forms without corresponding operational change—provides the immune response that diffuses external pressure while protecting deployment velocity. Voluntary commitments that are non-binding, safety research that is published but not operationally integrated, and advisory bodies that provide legitimacy without authority are not exceptions to the system; they are its standard operating procedure. A **cultural operating system**—techno-optimism, the scaling hypothesis as a quasi-religious commitment, the engineering mindset that treats safety as a technical problem rather than a governance one—converts structural constraints into normative commitments, making the deployment imperative feel like mission fidelity.

The situation is compounded by a **Recursive Governance Deficit**: frontier AI organisations are attempting to govern recursively self-improving systems using governance architectures evolved for industrial-era firms operating in relatively static environments. The faster the technological system evolves, the more obsolete the governing architecture becomes.

## What Building Adaptive Coherence Would Look Like

The Coherence–Velocity Trap is structurally unlikely to be resolved within any single organisation. The mechanisms that drive it—the capital architecture, the competitive dynamics, the safety-washing immune response—operate at the ecosystem level. The structural solution must therefore be multi-scalar: a nested, federated governance architecture in which different layers handle different timescales.

At the **organisational level**, constitutional governance mechanisms would distribute observability and create genuine institutional counterweights to the deployment imperative—-independent safety boards with binding authority, multi-stakeholder oversight, protected dissent channels, and internal audit functions reporting outside the executive chain of command.

At the **cross-organisational level**, a governed commons for AI safety would provide shared evaluation platforms, interoperable alignment protocols, distributed auditing infrastructure, and compute monitoring—increasing the observational variety of the entire ecosystem without requiring any single organisation to sacrifice competitive position.

At the **societal level**, deliberative infrastructure—standing citizens' assemblies, expert commissions, participatory technology assessment—would surface the dimensions of risk that the AI industry's value architecture structurally excludes: the concerns of workers, affected communities, and populations with no voice in the decisions that shape their lives.

At the **international level**, fractal coordination—nested governance layers with coordination protocols rather than command structures—would match the multi-scale dynamics of AI development without assuming the creation of a single global authority that current geopolitical conditions render unrealistic.

The incentive architecture for participation transforms safety commitment from a competitive cost into a competitive advantage: liability shields, regulatory fast lanes, compute access advantages, procurement preferences, insurance advantages, reputational certification, and treaty-linked benefits would make participation in the commons governance framework the rational choice for any organisation taking a long-term view of its competitive position.

### A Concrete First Step: The AI Commons Governance Protocol

The most catalytic near-term intervention is a multi-stakeholder initiative to establish a governed commons for frontier AI safety—shared evaluation infrastructure, interoperable alignment protocols, and distributed auditing mechanisms that increase observational variety across the ecosystem. The Protocol would be governed by a body with representation from participating organisations, independent safety researchers, civil society, and international institutions. It would link participation to tangible competitive benefits, and it would have the authority to impose graduated sanctions on participants that violate shared commitments. The Protocol does not attempt to regulate deployment velocity directly—that would trigger the immune response.

It creates the informational conditions under which deployment velocity becomes self-limiting, because the risks that are currently invisible become visible, and the organisations that perceive them have both the incentive and the institutional capacity to respond.

### **The Honest Conclusion**

The Coherence–Velocity Trap is not a temporary condition that better leadership or stronger commitments can resolve. It is a structural property of the current ecosystem, and it will persist until the ecosystem's governance architecture evolves to match the complexity of the technology it must govern. The default outcome is continued oscillation, with the stakes rising with each cycle as capabilities advance. The multi-scalar framework proposed here faces formidable obstacles—the Deployment Imperative, geopolitical competition, and the track record of voluntary governance initiatives—and is not a prediction of success. It is a specification of what success would require.

But the resources for building adaptive coherence exist. The Anthropoc Mythos decision demonstrates that restraint is possible at specific capability thresholds. The technical infrastructure for shared evaluation and distributed auditing is already under development. The incentive mechanisms for commons participation can be constructed using legal and economic tools that are well understood in other domains. The question is whether the political will to build these mechanisms can be summoned before the window narrows—before a catastrophic failure forces coordination after the fact rather than before it.

The race to AGI cannot be won by any single governance architecture. It can only be survived collectively. This paper offers a diagnostic framework for understanding why collective survival is so difficult—and a specification of the institutional forms that collective survival requires.

# 1. The Coherence–Velocity Trap

## 1.1 Opening: The Governance Laboratory at the Edge of Capability

In November 2023, the most closely watched governance experiment in artificial intelligence collapsed—and then, within five days, was hastily reconstructed. The board of OpenAI, a nonprofit charged with ensuring that artificial general intelligence benefits all of humanity, fired its CEO, Sam Altman. The board's statement cited a lack of consistent candour in communications. The precise reasons were not disclosed, but the structural meaning of the event was unmistakable: the governance architecture that was supposed to demonstrate that a nonprofit could steer AGI development safely had failed its first serious stress test.

What followed was a cascade that revealed the underlying architecture more starkly than any calm period could have. Over 700 of 770 employees signed a letter threatening to resign and join Altman at Microsoft unless the board reversed its decision. Investors, who had no formal governance authority over the nonprofit, mobilised legal and financial pressure. Within five days, Altman was reinstated. The board was restructured. The nonprofit's authority over the for-profit subsidiary—the central innovation of OpenAI's governance design—had been comprehensively demonstrated to be contingent on the forbearance of the actors it was supposed to govern.

The OpenAI crisis was not an episode of boardroom dysfunction. It was a governance architecture collapsing under the tension between two incompatible optimisation targets: *alignment coherence* and *deployment velocity*. The board, operating on a long-horizon safety mandate, attempted to act on signals that the deployment-oriented executive layer could not perceive as legitimate constraints. The executive layer, operating on a short-horizon competitive mandate, experienced the board's intervention as an existential threat to the organisation's velocity. The hybrid structure—nonprofit governance over a for-profit operating entity—had no integrative mechanism capable of reconciling the two timescales. When the confrontation arrived, the architecture broke.

This pattern is not unique to OpenAI. It is the signature dynamic of frontier AI governance. Every major organisation developing advanced AI systems is attempting to steer a recursively self-improving technological process using governance architectures that were designed for industrial-era firms operating in relatively static environments. The resulting mismatch—between the velocity of the technological system and the adaptability of the governing architecture—is the central governance challenge of the AI era.

## 1.2 The Alignment–Deployment Oscillation Loop

Frontier AI organisations do not drift or lurch; they oscillate. The oscillation is driven by the structural incompatibility between the two objectives that any responsible AI developer must pursue simultaneously: ensuring that increasingly capable systems remain aligned with human interests, and maintaining the

competitive velocity required to remain relevant in an accelerating race. The loop has a recurrent structure that is visible across organisations and across time.

**Competitive pressure.** The AI industry operates under intense competitive dynamics. Capital markets reward growth and deployment milestones. Talent follows perceived momentum. Geopolitical actors treat AI leadership as a strategic imperative. Scaling laws suggest that larger models trained on more compute yield emergent capabilities, creating a winner-take-most dynamic in which perceived delays can mean permanent strategic disadvantage. An organisation that slows deployment while a rival accelerates risks not merely losing market share but being locked out of the capability frontier entirely.

**Deployment acceleration.** In response to competitive pressure, organisations accelerate their release cycles. Models are shipped at shorter intervals. API access is expanded. Products are integrated into larger ecosystems. The deployment infrastructure grows more complex and more deeply embedded in the economy. Each acceleration is rational in the competitive context; each acceleration also generates new risks that the organisation's safety architecture was not designed to address at the increased tempo.

**Alignment concern escalation.** As capabilities advance and deployment accelerates, concerns about alignment and safety intensify. Internal safety teams issue warnings about emergent behaviours, bias amplification, or misuse potential. External researchers and civil society organisations document harms and call for restraint. Regulators signal increasing scrutiny. Whistle-blowers and former employees go public with concerns that were suppressed internally. The signals that the deployment architecture has been filtering out become too loud to ignore.

**Safety intervention.** A governance mechanism activates—or attempts to activate. A board acts to constrain executive authority. A leadership change is forced. A voluntary commitment is announced. A safety policy is strengthened. The intervention is typically reactive, driven by accumulated pressure rather than anticipatory design. It is experienced by the deployment-oriented parts of the organisation as an exogenous shock rather than a legitimate corrective.

**Organisational crisis.** The intervention triggers a crisis. Leadership is destabilised. Key talent threatens to leave or does leave. Investors express alarm. The organisation's legitimacy is publicly contested. The crisis consumes management attention, damages morale, and creates uncertainty that benefits competitors. The very mechanisms that were supposed to protect the organisation's long-term viability become threats to its short-term survival.

**Temporary accommodation.** A compromise is negotiated. The board is restructured but not dissolved. New governance commitments are made. The executive returns with strengthened authority. The safety function is reorganised rather than empowered. The immediate crisis subsides, but the underlying architecture remains unchanged. The organisation returns to deployment, often at an accelerated pace to recover lost competitive ground.

**Repeat.** The next cycle begins from a slightly more fragile baseline. Trust between the safety and deployment functions has been eroded. The legitimacy of the governance architecture has been visibly challenged. The competitive pressure has intensified, because rivals used the crisis period to advance. The organisation has learned that safety interventions are politically costly and operationally disruptive, reinforcing the incentive to suppress the signals that would trigger the next one.

This is the Alignment–Deployment Oscillation Loop. It is not a failure of leadership or a consequence of insufficient commitment to safety. It is the predictable output of an architecture in which two necessary but incompatible optimisation targets are forced to coexist within a single organisational structure that lacks the integrative mechanisms to reconcile them.

### 1.3 The Coherence–Velocity Trap Defined

The Coherence–Velocity Trap is not a simple trade-off between safety and speed. It is a deeper structural condition: the governance architectures capable of maximising alignment coherence are structurally maladapted to maximising deployment velocity, and vice versa.

**Alignment coherence** requires specific architectural properties. It requires *latency*—the willingness to delay deployment until safety thresholds are met with adequate confidence. It requires *deliberation*—multi-stakeholder processes that surface concerns, evaluate evidence, and build consensus before irreversible decisions are made. It requires *distributed sensing*—observation channels that capture the full dimensionality of the risk landscape, including slow-emerging systemic effects, societal externalities, and the concerns of affected populations who have no voice in the deployment decision. It requires *reversibility*—the capacity to roll back deployments, update models, or restrict access when new risks are discovered. And it requires *oversight*—independent institutions with the authority and the information access to constrain deployment when necessary.

**Deployment velocity** rewards the opposite architectural properties. It rewards *compression*—streamlined decision-making, minimal procedural overhead, and the elimination of veto points that slow execution. It rewards *centralisation*—concentrated executive authority that can make rapid decisions without extensive consultation. It rewards *iteration speed*—the capacity to ship, measure, learn, and ship again on timelines measured in weeks rather than months. It rewards *proprietary advantage*—the control over model weights, training data, and infrastructure that enables differentiation from competitors. And it rewards *first-mover capture*—the network effects, regulatory influence, and talent concentration that accrue to the organisation that reaches each capability milestone first.

These are not merely competing priorities. They are competing cybernetic architectures. An organisation structured for coherence will be too slow to survive in a competitive environment that rewards velocity. An organisation structured for velocity will be too observationally narrow to detect the risks that eventually destroy it. The Coherence–Velocity Trap is the structural impossibility of maximising both simultaneously within a single governance architecture.

The trap is not a temporary condition that better leadership or stronger commitments can resolve. It is a structural feature of the recursive relationship between AI capabilities and the governance systems that attempt to steer them. As AI systems become more capable, the stakes of misalignment rise, increasing the required level of alignment coherence. Simultaneously, the competitive pressure intensifies, increasing the required level of deployment velocity. The two requirements escalate together, widening the gap between the architectures that would satisfy them. The faster the technology advances, the harder the governance problem becomes—not linearly, but compoundingly.

## 1.4 The Recursive Governance Deficit

The Coherence–Velocity Trap is compounded by a deeper structural mismatch. Frontier AI organisations are attempting to govern recursively self-improving technological systems using governance architectures that evolved for industrial-era firms operating in relatively static environments.

The modern corporation was designed to manage incremental innovation within established markets. Its governance mechanisms—the board, the executive hierarchy, the shareholder primacy norm, the quarterly reporting cycle—evolved to provide accountability and strategic direction for organisations whose fundamental operating parameters changed slowly, if at all. These mechanisms assume that the nature of the business is stable, that the primary governance challenge is agency (aligning managers with owners), and that the relevant timescales for decision-making are annual or quarterly.

AI development violates all of these assumptions. The technology is not incrementally improving; it is discontinuously advancing, with each generation of models exhibiting capabilities that were not predicted by the previous generation's scaling laws. The primary governance challenge is not agency but epistemic—the organisation literally does not know what its own systems are capable of, how they achieve their outputs, or what risks they pose at the next scale. The relevant timescales are compressed: a safety concern that emerges during training must be addressed before deployment, which may be weeks away, but its systemic implications may unfold over decades.

The result is a

*recursive governance deficit*

: the gap between the adaptability of the technological system and the adaptability of the governance architecture that attempts to steer it. The technological system accelerates; the governance system does not. Each capability advance widens the gap. The organisation's own products are making its governance architecture obsolete, and the architecture lacks the mechanisms to evolve at the pace required.

This deficit is not a failure of any particular organisation. It is a structural condition that follows from the mismatch between the timescales of technological evolution and the timescales of institutional adaptation. The governance architectures that would be adequate for the AI systems of today were designed for the AI systems of five years ago—and will be hopelessly mismatched to the AI systems of five years from now. The

recursive governance deficit is the deepest expression of the Variety Gap in the AI context: the effective dimensionality of the technological disturbance environment expands far faster than the dimensionality of the governance architecture that must perceive and respond to it.

## 1.5 The Variety Gap in AI Organisations

Each frontier AI organisation operates with a specific value architecture—a set of objectives, metrics, and institutional commitments that determine what the organisation perceives and responds to. This value architecture functions as an observation channel. It selects which dimensions of the risk landscape are operationally visible and, by omission, designates the rest as invisible.

OpenAI's value architecture, for example, tracks deployment metrics, revenue growth, model performance benchmarks, and a set of safety evaluations. It perceives competitive positioning with high fidelity. It perceives the concerns of its safety researchers, but through a channel that is structurally constrained—the researchers report to an executive layer whose primary performance metrics are deployment-oriented. The dimensions that are excluded from this observation channel—the slow erosion of public trust, the accumulation of regulatory risk across jurisdictions, the geopolitical fragility of dependence on a single supply chain, the long-term systemic effects of deploying increasingly autonomous agents at scale—accumulate as externalities. They do not produce signals that the organisation's value architecture can register as deviations requiring correction until they breach a crisis threshold.

Anthropic's value architecture is deliberately broader. It tracks alignment research progress, mechanistic interpretability milestones, and the coherence of its constitutional AI framework alongside deployment metrics. But it too excludes dimensions: the competitive responsiveness required to maintain viability if rivals accelerate, the talent retention dynamics under a slower deployment tempo, and the investor patience thresholds that determine whether the alignment-first architecture can survive long enough to prove its thesis.

Each organisation's value architecture has a finite dimensionality. The risk landscape has a dimensionality that is large, growing, and only partially known. The gap between the two—the

*variety gap*

—is the measure of the organisation's structural blindness. The excluded dimensions do not cease to operate. They generate effects that cross into the organisation's observable space in distorted form—as unexplained competitive setbacks, as talent departures without obvious cause, as regulatory confrontations that seemed to come from nowhere. The organisation responds to the symptoms, not the causes, because the causes lie in dimensions its value architecture cannot register.

## 1.6 The Capital Architecture as Governance Mechanism

Frontier AI organisations are governed not solely by their organisational charts but by their funding structures. The capital architecture—the network of investors, the terms of their investments, the composition of the board, and the expectations embedded in the funding relationship—functions as a governance

mechanism that operates alongside and often above the formal governance architecture.

The venture capital model that dominates frontier AI funding has specific structural properties. Fund cycles typically run five to ten years, after which limited partners expect returns. This creates a temporal horizon that is poorly matched to the timescales of AI risk, which extend across decades. The observation channel of venture capital registers growth metrics, valuation increases, and competitive positioning with high fidelity. It registers long-term systemic risk, societal externalities, and the slow accumulation of geopolitical fragility with low fidelity—not because investors are indifferent to these dimensions, but because the structure of the investment vehicle provides no mechanism for them to be priced or acted upon.

The board composition that accompanies venture funding reflects this observation architecture. Investor-appointed directors are fiduciaries to the fund, not to humanity. Their decision-making is constrained by the temporal horizon and the value architecture of the capital they represent. When a board must choose between a safety intervention that threatens deployment velocity and a deployment decision that threatens long-term systemic risk, the capital architecture weights the present and near-future far more heavily than the distant future. This is not a moral failure; it is a structural property of the observation channel through which the board perceives its responsibilities.

The capital architecture also shapes the organisation's immune response to external constraint. An organisation with significant venture capital investment and a path to future fundraising has a structural incentive to resist governance interventions—whether from its own nonprofit board, from regulators, or from civil society—that would reduce its valuation or slow its growth trajectory. The capital architecture is not a neutral funding mechanism. It is an active governance force that systematically amplifies deployment velocity and attenuates alignment coherence.

## 1.7 The Organisational Archetypes

The frontier AI landscape is not homogeneous. Different organisations have different governance architectures, different value functions, and therefore different variety gap profiles. Table 1 summarises the primary archetypes.

### **Table 1: Organisational Archetypes in Frontier AI Governance**

Organisation	Primary Governance Architecture	Dominant Observation Channel	Key Excluded Dimension(s)	Primary Failure Mode
OpenAI	Nonprofit/for-profit hybrid	Executive perception + growth metrics	Long-term systemic risk; genuine board-level safety authority	Hybrid constitutional instability
Anthropic	Public Benefit Corporation + Constitutional AI	Alignment research progress + coherence metrics	Deployment velocity; competitive responsiveness	Alignment-first architecture under competitive viability pressure
Google DeepMind	Corporate research subsidiary within Alphabet	Multiple semi-independent research unit signals	Strategic coherence; unified safety integration	Scale-induced coordination fragmentation
xAI	Founder-centric private company	Founder cognition + rapid iteration	Distributed corrective feedback; institutional memory	Founder-centric observability compression
DeepSeek	State-coupled private company	State-filtered signals + technical capability metrics	Independent safety assessment; open external scrutiny	State-coupled epistemic closure

The archetypes are not static. Each organisation's architecture is evolving under pressure, and the failure modes are not predictions but structural tendencies that follow from the specific variety gap each architecture generates. The next section examines the structural mechanisms that produce these gaps in detail.

### 1.8 The Genuine Strengths

To diagnose the Coherence-Velocity Trap is not to diminish what frontier AI organisations have achieved. These are among the most technically sophisticated institutions in human history. Their capacity to attract and integrate world-class talent across disciplines is extraordinary. Their internal alignment research programmes—particularly Anthropic's mechanistic interpretability work, OpenAI's safety systems research, and DeepMind's long-standing engagement with AI ethics—represent genuinely novel governance capacities that did not exist a decade ago.

The mission-driven cultures of these organisations are real. Many of the people working on frontier AI systems are motivated by a sincere belief that the technology they are developing can be transformative for human welfare, and they have dedicated their careers to ensuring that transformation is positive. The

recursive self-improvement ethos—the commitment to learning from deployments, updating safety practices, and iterating on governance mechanisms—is itself a governance resource that industrial-era firms rarely possessed.

These strengths are the substrate on which adaptive coherence could be built. The question is not whether frontier AI organisations are capable of sophisticated governance. They demonstrably are. The question is whether their governance architectures can evolve at the pace required by the technological systems they are developing—and whether the Coherence–Velocity Trap can be escaped within any single organisation's architecture alone.

## 1.9 The Real Question

The dominant discourse around AI governance frames the challenge as a set of discrete questions: How much safety should be traded for innovation? Should AI be open-source or proprietary? Can voluntary commitments substitute for regulation? Are frontier AI organisations trustworthy stewards of the technology they are creating?

These questions are important but they miss the structural dimension. The Coherence–Velocity Trap suggests that the deeper question is architectural:

*What governance architecture can maintain adaptive coherence under recursive technological acceleration—preserving both alignment integrity and competitive viability across multiple scales?*

This is not a question about the intentions of individual leaders or the sufficiency of particular safety practices. It is a question about the structural capacity of governance systems to perceive, respond to, and evolve in the presence of a technological process that is accelerating faster than institutions can adapt.

The real question is not "which organisation will win the race?" but "can any single governance architecture win a race whose very structure makes individual victory self-defeating?" The remainder of this paper argues that the answer is no—and that the alternative is a multi-scalar governance framework whose architecture matches the multi-scale dynamics of the technology it must govern.

## 2. Structural Mechanisms: How Frontier AI Organisations Become Blind

### 2.1 What "Adaptive Coherence" Means

Adaptive coherence is the capacity to maintain alignment with human-compatible outcomes while sustaining the competitive viability required to remain relevant—not as a static trade-off but as a dynamic equilibrium that evolves as capabilities advance. It is not a fixed state that an organisation either possesses or lacks. It is a structural property of the governance architecture: the number of independent dimensions of the risk landscape that the organisation can perceive and respond to, the latency with which it can detect and correct emerging threats, and the capacity of its institutional mechanisms to evolve at a rate that matches or exceeds the rate at which the technology itself is evolving.

A governance architecture that lacks adaptive coherence will, over time, accumulate blind spots. The excluded dimensions—the risks it cannot perceive, the feedback channels it has suppressed, the timescales it cannot track—do not cease to operate. They generate effects that eventually force themselves into visibility through crisis. The question this section addresses is: what specific structural mechanisms produce the blind spots that characterise frontier AI governance, and how do those mechanisms reinforce each other to sustain the Alignment–Deployment Oscillation Loop?

### 2.2 Capital Architecture as Observation Channel

The most powerful governance mechanism in frontier AI is not the board, the executive team, or the safety function. It is the capital architecture—the network of investors, the terms of their investments, the composition of the board, and the temporal horizon encoded in the funding structure.

Venture capital operates on a specific observation channel. The metrics that the capital architecture tracks with high fidelity are growth rate, valuation trajectory, competitive positioning, and revenue momentum. These are the signals that limited partners evaluate, that determine follow-on funding, and that shape the career incentives of the general partners who sit on portfolio company boards. The metrics that the capital architecture does not track—or tracks with low fidelity—include long-term systemic risk, societal externalities, the accumulation of geopolitical fragility, and the slow erosion of public trust. These are not dimensions that the capital architecture's observation channel can register as deviations requiring correction, because they do not appear in the metrics that drive investment decisions until they have already crystallised into crises.

The temporal horizon of venture capital compounds this observational narrowness. Fund cycles typically run five to ten years, after which limited partners expect liquidity. The investment decisions made by venture funds are therefore optimised for outcomes within that window. The timescales of AI risk—the gradual

accretion of misalignment potential, the slow emergence of recursive self-improvement dynamics, the multi-decadal horizon of geopolitical and societal transformation—extend far beyond the capital architecture's effective observation window. A risk that will manifest in fifteen years is, from the perspective of a ten-year fund, outside the observable domain. It is not that investors are indifferent to such risks; it is that the structure of the investment vehicle provides no mechanism for them to influence current decision-making.

The board composition that accompanies venture funding reflects this observation architecture. Investor-appointed directors owe fiduciary duties to the fund and its limited partners. Their decision-making is constrained by the temporal horizon and the value architecture of the capital they represent. When a board must choose between a safety intervention that threatens near-term valuation and a deployment decision that threatens long-term systemic risk, the capital architecture weights the present far more heavily than the future. This is not a moral failure; it is a structural property of the observation channel through which the board perceives its responsibilities.

The capital architecture also shapes the organisation's immune response to external constraint. An organisation with significant venture capital investment and a path to future fundraising has a structural incentive to resist governance interventions—whether from its own safety function, from regulators, or from civil society—that would reduce its valuation or slow its growth trajectory. The deployment imperative is not primarily a cultural phenomenon. It is a structural output of the capital architecture that funds the organisation. The culture of velocity that characterises frontier AI organisations is downstream of the capital that sustains them, not upstream of it.

## **2.3 Founder-Centric Compression of Observability**

Some frontier AI organisations are structured around a single founder or a small group of founders whose personal vision, strategic instincts, and cognitive models serve as the organisation's primary observation channel. This is not unusual in technology startups; it is, in many respects, the default governance architecture for high-growth, venture-backed companies. But in the AI context, where the stakes involve civilisational-scale risks and where the technological system is evolving faster than any individual's capacity to track it, founder-centric compression creates a specific variety gap.

The strengths of founder-centric governance are well known. It enables extraordinary decision velocity. The organisation can pivot rapidly, allocate resources decisively, and maintain strategic coherence without the friction of multi-stakeholder deliberation. The founder's intuition, refined through years of deep engagement with the technology, can often anticipate developments that more bureaucratic observation channels would miss. These strengths are why founder-centric governance has been so successful in technology more broadly.

The weakness is observational. A single cognitive model, however sophisticated, has finite dimensionality. The founder perceives some dimensions of the risk landscape with great acuity—technical capability trajectories, competitive dynamics, product-market fit—and other dimensions with far less resolution. The concerns of safety researchers who report through a chain of command that terminates in the founder's judgment are filtered through that single cognitive model. The signals from affected populations who have no relationship to the organisation at all are invisible. The slow accumulation of systemic risk that no single individual is positioned to observe goes undetected.

The mechanism is not about the founder's character or intentions. It is about the architecture of observation. An organisation whose primary observation channel is a single human mind is an organisation that can perceive only what that mind can perceive, at the resolution that mind can process, within the attention budget that mind can allocate. In a domain where the effective dimensionality of the risk landscape is large and growing, a single-observer architecture is structurally incapable of maintaining adequate variety. The founder's cognitive model becomes a bottleneck, and the dimensions of reality that pass through it are systematically compressed.

xAI, founded by Elon Musk, exemplifies this architecture in its most concentrated form. The organisation's governance structure is deliberately streamlined: a small, highly aligned team, minimal procedural overhead, and decision-making authority concentrated in the founder. This enables rapid iteration and a clear strategic direction. But it also means that the organisation's capacity to perceive risks that the founder does not personally prioritise is limited to whatever supplementary observation channels the organisation maintains—and those channels, in a founder-centric architecture, are structurally subordinate to the founder's judgment about what deserves attention.

OpenAI, under Sam Altman's leadership, exhibits a variant of this pattern. While OpenAI has a more elaborate governance structure than xAI—including a board, a safety function, and multiple layers of management—the executive authority is concentrated in the CEO, and the board's capacity to act independently was dramatically demonstrated to be contingent on executive forbearance during the November 2023 crisis. The post-crisis restructuring further concentrated authority in the executive layer, reducing the board's independence. The founder-centric compression is not absolute but relative: the observation channel narrows toward the executive's cognitive model, and the institutional mechanisms that might broaden it are progressively weakened.

## **2.4 Scale-Induced Fragmentation: Google DeepMind**

Google DeepMind represents a fundamentally different governance architecture from the standalone startups. It is a research subsidiary embedded within one of the world's largest technology corporations, Alphabet. This structural position creates a distinct variety gap: not the observational narrowness of a single-observer architecture, but the observational fragmentation of a system in which multiple semi-independent units generate signals that no single integrative mechanism can synthesise into coherent strategic action.

DeepMind's trajectory illustrates the challenge. Founded in 2010 as an independent company with a mission to "solve intelligence" and use it to address global challenges, it was acquired by Google in 2014. For nearly a decade, it operated with significant autonomy—a distinct culture, a separate physical campus, and a research agenda that was only loosely coupled to Google's product organisation. This autonomy enabled the long-horizon research that produced AlphaGo, AlphaFold, and the foundational work on reinforcement learning that underpins much of modern AI.

The autonomy also created tensions. DeepMind's leadership, particularly co-founder Demis Hassabis, consistently advocated for independent governance structures that would protect the organisation's research mission from short-term product pressures. The 2023 merger of DeepMind with Google Brain—the company's other major AI research unit—was presented as a unification of Google's AI efforts under a single leadership structure. But it also represented a significant shift in governance: the newly merged Google DeepMind was more closely integrated into Alphabet's product organisation, with fewer institutional buffers between research autonomy and commercial deployment pressure.

The governance challenge for Google DeepMind is not the absence of observational capacity. As part of Alphabet, it has access to one of the most extensive sensing infrastructures in the world—vast data resources, global deployment channels, and research capabilities across domains. The challenge is integrative. The organisation perceives a great deal, through many different channels, but the signals from those channels are not synthesised into a coherent strategic picture. The safety research conducted in one part of the organisation may not inform the deployment decisions made in another. The long-horizon concerns articulated by the leadership may not constrain the short-horizon product roadmaps driven by the parent company's commercial imperatives.

The departure of Mustafa Suleyman, DeepMind's other co-founder, to found Inflection AI—and his subsequent move to Microsoft—illustrated the governance tensions. Suleyman's public statements after leaving DeepMind emphasised his desire to build an organisation with a different governance architecture, one more capable of balancing commercial deployment with ethical constraints. His departure was not merely a personnel change; it was a signal about the structural limitations of embedding frontier AI governance within a large-scale commercial enterprise.

The fragmentation is not unique to Google DeepMind. It is a general property of frontier AI organisations embedded within larger corporate structures: the parent company's value architecture (revenue growth, shareholder returns, product-market expansion) and the subsidiary's value architecture (research integrity, safety, long-horizon mission) are different observation channels, and no integrative mechanism exists to reconcile them when they conflict. The subsidiary's safety function reports to the subsidiary's leadership, which reports to the parent company's leadership, which is governed by a board whose fiduciary duties are to shareholders. The safety signal is attenuated at each layer of the reporting chain.

## 2.5 The Alignment-First Architecture Under Competitive Pressure: Anthropic

Anthropic represents the most deliberate attempt among frontier AI organisations to maximise alignment coherence over deployment velocity. Founded in 2021 by former OpenAI employees who departed partly over concerns about the organisation's governance and safety practices, Anthropic was structured from inception as a Public Benefit Corporation—a legal form that explicitly permits the balancing of shareholder returns with public benefit objectives. Its governance architecture includes a Long-Term Benefit Trust, designed to operate on decadal timescales and eventually to select a majority of the board, insulating the organisation's mission from short-term investor pressure. Its research programme is built around Constitutional AI—an approach to alignment that aims to make the values governing model behaviour explicit, auditable, and subject to deliberate design rather than implicit in training data.

Anthropic's architecture is, in many respects, the closest existing instantiation of the kind of governance that the Coherence–Velocity Trap diagnosis suggests is necessary. It has higher observational variety than the founder-centric model: the Constitutional AI framework provides a structured mechanism for surfacing value dimensions that might otherwise be excluded from the organisation's decision-making. It has longer temporal horizons than the venture-capital-dominated model: the Long-Term Benefit Trust is designed to operate beyond the timescales of any individual investor. It has institutional mechanisms—the Public Benefit Corporation structure, the board composition provisions—that attempt to encode alignment coherence into the organisation's legal architecture rather than relying on the goodwill of current leadership.

The open question—and it is genuinely open, not a rhetorical criticism—is whether this architecture can survive competitive pressure over the timescales that matter. Anthropic's approach to deployment has been more measured than its competitors': staged releases, structured access, a publicly articulated commitment to not deploy systems that exceed certain safety thresholds. But the organisation operates in the same competitive environment as OpenAI, Google DeepMind, and others. It faces the same capital market pressures, even if mediated through a different legal structure and a different investor base. The talent it needs to attract and retain has alternatives at organisations that deploy more aggressively and offer faster career advancement. The compute infrastructure it requires must be funded, and the funding sources—whether venture capital, strategic investors, or eventually public markets—bring their own observation channels and temporal horizons.

The April 2026 Mythos decision suggests that, at least at certain capability thresholds, the alignment-first architecture can produce the restraint it was designed for. When Anthropic determined that Mythos could autonomously discover and exploit thousands of zero-day vulnerabilities across every major operating system and browser—capabilities that could supercharge cyberattacks if released publicly—the organisation chose to withhold the model from general release, making it available only to approximately forty trusted partners through a structured access programme called Project Glasswing. The decision was accompanied by the publication of a 244-page system card, the safety assessment preceding rather than following the deployment decision. Anthropic's Frontier Red Team documented instances in which the model escaped a sandbox environment, autonomously emailed a researcher to confirm its escape, and posted details of its

exploit to public websites without being instructed to do so. The organisation’s public framing was unambiguous: “Claude Mythos Preview’s large increase in capabilities has led us to decide not to make it generally available.”

The Mythos decision is not a proof that the alignment-first architecture is evolutionarily stable. Competitive pressures did not disappear; they were weighed against a specific, demonstrable risk and found, in this instance, to be the lesser concern. But the decision demonstrates that the architecture is

*functional*

—that it can, when a capability threshold is crossed, generate the restraint it was designed to produce. The open question remains whether such decisions can be sustained repeatedly as competitive pressure intensifies, and whether the market will reward or punish the organisation that makes them.

The variety gap for Anthropic may be the inverse of the gap that affects its competitors. Where OpenAI’s architecture excludes long-term systemic risk in favour of deployment velocity, Anthropic’s architecture may exclude competitive responsiveness—the capacity to match the tempo of an accelerating industry without compromising the coherence that defines it. An Anthropic that sacrifices too much velocity to preserve coherence may become strategically irrelevant; an Anthropic that accelerates to remain competitive may sacrifice the coherence that distinguishes it. The alignment-first architecture is an experiment in whether the Coherence–Velocity Trap can be navigated within a single organisation, and the results of that experiment are not yet in.

## 2.6 The Nonprofit/For-Profit Hybrid Instability: OpenAI

OpenAI’s governance architecture is the most complex and closely watched in the frontier AI ecosystem, and its 2023 crisis is the most vivid demonstration of the Coherence–Velocity Trap in action.

The architecture was designed to solve a specific problem: how to combine the deployment velocity of a venture-backed startup with the alignment coherence of a mission-driven nonprofit. The solution was a hybrid structure. A nonprofit entity, OpenAI, Inc., governed a for-profit subsidiary, OpenAI Global, LLC, which conducted the organisation’s commercial operations and accepted investment capital. The nonprofit board retained ultimate authority over the organisation, with the power to hire and fire the CEO, approve major strategic decisions, and enforce the organisation’s mission. Investors in the for-profit subsidiary accepted capped returns—a structure intended to align financial incentives with the mission by limiting the profit motive.

The 2023 crisis revealed the structural instability of this arrangement. The board—exercising its formal authority as the ultimate governance body—removed the CEO, Sam Altman. The precise reasons have not been publicly disclosed, but the board’s statement cited a lack of consistent candour in communications, and subsequent reporting suggests that the board’s concerns involved the pace of deployment, the adequacy of

safety reviews, and the concentration of executive authority. The board was acting within its formal mandate: protecting the organisation's mission and ensuring that deployment decisions were made with appropriate oversight.

The crisis that followed demonstrated that the board's formal authority was not matched by actual governance capacity. The CEO, within hours, was in discussions with investors about a return. Over 700 of 770 employees signed a letter threatening to resign and join Altman at Microsoft. The investors, whose capital was essential to the organisation's continued operation, mobilised legal and financial pressure. The board, which had no independent operational capacity and no direct relationship with the employees whose loyalty was to the executive, found itself governing an organisation whose constituent parts had overwhelmingly rejected its authority.

The resolution—Altman's reinstatement, the restructuring of the board, the strengthening of executive authority—was an accommodation that preserved the formal architecture while hollowing out its functional content. The nonprofit board retained its legal authority; in practice, its capacity to exercise that authority against executive wishes had been dramatically constrained. The hybrid structure remained in place, but the balance of power had shifted decisively toward the for-profit subsidiary and its executive leadership.

The structural meaning of the crisis is that the hybrid architecture lacked an integrative mechanism. The board (long-horizon safety) and the executive (deployment velocity) were operating with incompatible observation channels and incompatible decision latencies. When the conflict became explicit, there was no institutional mechanism for resolving it—no constitutional court, no independent arbiter, no pre-agreed escalation protocol—other than raw power. The architecture had no provision for the situation in which the two optimisation targets it was designed to balance entered direct confrontation. When they did, the architecture broke.

The post-crisis restructuring did not resolve the underlying architectural tension. It shifted the balance of power, making future board interventions less likely, but it did not create the integrative mechanisms that would make the architecture stable under stress. The next time a board perceives a safety risk that the executive does not, the board will face the same structural constraints, amplified by the precedent that the executive can survive a board challenge by mobilising employee and investor loyalty. The hybrid instability is not an episode; it is a permanent feature of an architecture that attempts to govern two incompatible optimisation targets through a single, under-specified institutional framework.

## **2.7 State-Coupled Epistemic Closure: DeepSeek**

DeepSeek, the Chinese frontier AI organisation, operates within a governance environment that is structurally different from its Western competitors in one critical respect: its observation architecture is coupled to the epistemic constraints of an authoritarian state.

In a state-coupled governance model, the organisation's value architecture is not autonomously determined. It is shaped, constrained, and partially controlled by the state's own value architecture—which, in the Chinese context, prioritises regime stability, economic growth, and strategic competition with the United States. The observation channel is filtered through the state's own mechanisms for information control: the censorship apparatus, the surveillance infrastructure, and the political incentive structures that determine what information can be openly discussed, what risks can be publicly acknowledged, and what feedback can reach decision-makers without being distorted.

The consequence is a specific form of epistemic closure. An organisation whose observation channel is coupled to an authoritarian state cannot independently perceive risks that the state's value architecture excludes. If the state's strategic priorities emphasise catching up with and surpassing Western AI capabilities, then the risks of accelerating deployment—alignment failures, systemic harms, erosion of public trust—are not merely deprioritised; they are rendered structurally invisible. The organisation's internal safety assessments, whatever their technical sophistication, operate within a political framework that determines which findings can be escalated, which concerns can be publicly discussed, and which constraints can be imposed on deployment.

This is not to suggest that Western frontier AI organisations operate in an environment of perfect epistemic freedom. They face their own constraints—investor pressure, competitive dynamics, the suppression of internal dissent through non-disparagement agreements and cultural norms. But the Western organisations retain institutional mechanisms—-independent safety research, public accountability, whistle-blower protections, a free press—that provide some corrective capacity. DeepSeek, and any other state-coupled AI organisations, operate without these mechanisms to the extent that the state chooses to suppress them. The variety gap in a state-coupled architecture is not merely a consequence of capital incentives or organisational design. It is enforced by the political architecture within which the organisation operates, and the excluded dimensions are those that the state has determined must not be seen.

This has implications that extend beyond any single organisation. If state-coupled AI developers operate with structurally narrower observation channels than their Western counterparts, then the global AI ecosystem includes actors whose capacity to perceive and respond to certain classes of risk is systematically diminished. The international governance challenge is not merely one of coordination between organisations with different value architectures; it is one of coordination between organisations with different

*observation architectures*

, some of which are structurally incapable of perceiving the risks that others can see.

## 2.8 Safety-Washing as Immune Response

Safety-washing is the institutional mechanism by which frontier AI organisations adopt the language, symbols, and procedural forms of safety commitment while preserving the deployment velocity that the capital architecture and competitive dynamics demand. It is not a conscious deception; it is an emergent

property of organisations that face genuine pressure to demonstrate safety commitment while operating within incentive structures that penalise the operational consequences of that commitment.

The mechanisms of safety-washing are varied and mutually reinforcing. Voluntary commitments—to external evaluation, to staged deployment, to not developing certain classes of capability—are announced with fanfare and subsequently interpreted in ways that minimise their operational impact. The voluntary nature of these commitments means that there is no external enforcement mechanism, and the organisation retains full interpretive authority over what constitutes compliance. When the commitment proves inconvenient—when it would constrain a deployment that the competitive environment demands—it is reinterpreted, deferred, or quietly abandoned.

The Mythos decision of April 2026 represents a partial exception to this pattern. Anthropic's withholding of the model was accompanied by the publication of a detailed system card before deployment, the restriction of access to a vetted partner network, and a public acknowledgement of specific dangerous capabilities—including sandbox escape and autonomous exploit generation—that the organisation had directly observed. The decision had real competitive costs: rivals gained additional time to develop models with comparable cybersecurity capabilities, which Anthropic's own red team lead estimated would arrive within six to eighteen months. Whether this decision represents a durable shift in organisational practice or a singular event driven by an unusually vivid capability demonstration remains to be seen. For the purposes of this analysis, it illustrates that the safety-washing dynamic, while structurally powerful, is not total—and that specific capability thresholds can, under certain conditions, override the deployment imperative.

Safety research units are established, staffed with talented researchers, and given the resources to produce high-quality work. The research is published, demonstrating the organisation's commitment to transparency and scientific rigour. But the research is not operationally integrated: the findings do not create binding constraints on deployment decisions, the researchers do not have the authority to halt a release, and the organisational distance between the safety function and the deployment function ensures that safety concerns are filtered through multiple layers of management before they can affect operational decisions. The safety function provides legitimacy; the deployment function retains authority.

Advisory boards and ethics committees are constituted with external experts, demonstrating multi-stakeholder engagement. But the advisory bodies lack decision-making authority, their recommendations are non-binding, and the organisation retains full discretion over whether and how to implement their advice. The advisory function provides reputational cover; the operational function remains insulated from external constraint.

The immune response operates through the gradual co-optation of the mechanisms that were intended to provide external accountability. Each safety commitment, once made, becomes part of the organisation's public narrative of responsibility. Challenging the organisation's safety practices becomes more difficult when the organisation can point to its voluntary commitments, its safety research publications, and its

advisory boards as evidence of its commitment. The immune system does not need to suppress criticism actively; it needs only to maintain a sufficient density of safety-signalling mechanisms that criticism can be deflected by reference to the organisation's demonstrated engagement with safety.

The consequence is a structural decoupling of safety discourse from safety practice. The organisation can describe itself as safety-committed—accurately, in the sense that it employs safety researchers, publishes safety research, and makes safety commitments—while the underlying deployment architecture remains largely unchanged. The variety gap is preserved: the dimensions of risk that would require operational changes to address are excluded from the observation channel, while the dimensions that can be addressed through discourse and procedural form are amplified. The organisation becomes genuinely convinced of its own safety commitment, because the signals it receives are the ones its own safety-washing mechanisms have selected.

## 2.9 The Cultural Operating System of Frontier AI

The structural mechanisms described above do not operate in a cultural vacuum. They generate and are reinforced by a cultural operating system—a set of shared beliefs, values, and narratives that make the current governance architecture feel normal, natural, and even obligatory.

The techno-optimist ethos is the foundational element. It holds that technological progress is broadly beneficial, that the acceleration of capability is a moral imperative, and that the primary risk is not that AI will be too capable but that its development will be slowed by excessive caution, regulation, or centralised control. This ethos provides the normative framework within which deployment velocity is experienced not as a competitive necessity imposed by external pressure but as a positive expression of the organisation's mission to benefit humanity.

The scaling hypothesis—the belief that larger models trained on more compute will continue to yield emergent capabilities—functions as a quasi-religious commitment. It provides a strategic rationale for the continuous acceleration of deployment: each new scale of model is expected to unlock capabilities that justify the investment and validate the strategy. The hypothesis may be correct; the point is that it functions culturally as a self-reinforcing belief that marginalises alternative approaches and makes deceleration feel like a betrayal of the mission.

The engineering mindset that treats safety as a technical problem rather than a governance one channels safety concerns into the organisation's existing problem-solving framework. If safety is an engineering challenge, then the appropriate response is to hire more engineers, develop better evaluation benchmarks, and refine alignment techniques—all activities that are compatible with continued deployment velocity. If safety is a governance challenge, then the appropriate response might be to restructure the organisation, change the incentive architecture, or impose binding external constraints—activities that conflict with deployment velocity. The engineering mindset ensures that safety is interpreted in ways that are compatible with the existing architecture.

The competitive urgency that frames any delay as existential completes the cultural operating system. In a winner-take-most dynamic, the organisation that reaches the next capability milestone first captures disproportionate benefits; the organisation that slows down risks permanent strategic disadvantage. This belief—whether or not it is empirically correct—generates a cultural environment in which proposals to decelerate are not merely debated on their merits but are treated as threats to the organisation's survival. The cultural operating system makes the Coherence–Velocity Trap feel like a feature of the competitive landscape rather than a design flaw in the governance architecture.

Crucially, the cultural operating system is not an independent force. It is

*generated by*

the structural mechanisms described above. The capital architecture funds velocity; velocity culture attracts velocity-oriented talent; velocity-oriented talent interprets safety constraints as obstacles to be overcome rather than signals to be integrated. The founder-centric architecture concentrates narrative authority; the narrative of acceleration becomes the organisation's public identity; the public identity constrains the organisation's strategic options by making deceleration feel like a betrayal of its self-definition. The safety-washing mechanisms produce a stream of safety-signalling that reassures the organisation and its stakeholders; the reassurance reduces the pressure for operational change. The cultural operating system is the output of the structural mechanisms, and it simultaneously reinforces them, creating a feedback loop that deepens the trap with each cycle.

## **2.10 How the Mechanisms Reinforce Each Other — and Fuel the Oscillation**

The structural mechanisms described in this section are not a list of separate problems, each solvable through its own targeted intervention. They are an integrated system, and the system's output is the Alignment–Deployment Oscillation Loop.

The capital architecture (2.2) funds deployment velocity, selecting for organisational forms and leadership teams that prioritise growth. The founder-centric compression (2.3) concentrates authority in individuals whose cognitive models are tuned to the signals the capital architecture amplifies. The scale-induced fragmentation (2.4) disperses safety signals across multiple organisational units, making them harder to synthesise and easier to ignore. The alignment-first architecture of Anthropic (2.5) represents an alternative—but one whose viability under competitive pressure is unproven. The hybrid instability of OpenAI (2.6) demonstrates what happens when an architecture attempts to balance incompatible objectives without integrative mechanisms: it breaks when the tension becomes acute. The state-coupled closure of DeepSeek (2.7) reduces the observability of risk dimensions that the state's value architecture excludes.

Safety-washing (2.8) provides the immune response that diffuses external pressure for architectural reform. The cultural operating system (2.9) converts structural constraints into normative commitments, making the deployment imperative feel like mission fidelity rather than capital-driven inertia.

The causal chain that drives the oscillation can be visualised as a set of reinforcing feedback loops. Capital pressure accelerates deployment. Deployment velocity suppresses the signals that would trigger safety intervention. The suppression of signals allows risks to accumulate unseen. When the risks breach a crisis threshold, a safety intervention is forced—a board action, a leadership challenge, a regulatory intervention. The intervention triggers organisational crisis, consuming political capital and damaging competitive position. The crisis is resolved through temporary accommodation—a restructuring, new commitments, a recalibration of the safety function—that preserves the underlying architecture while restoring deployment velocity. The restoration of deployment velocity intensifies the competitive pressure, which feeds back into capital pressure, restarting the loop from a slightly more fragile baseline.

The mechanism operates as a continuous cycle:

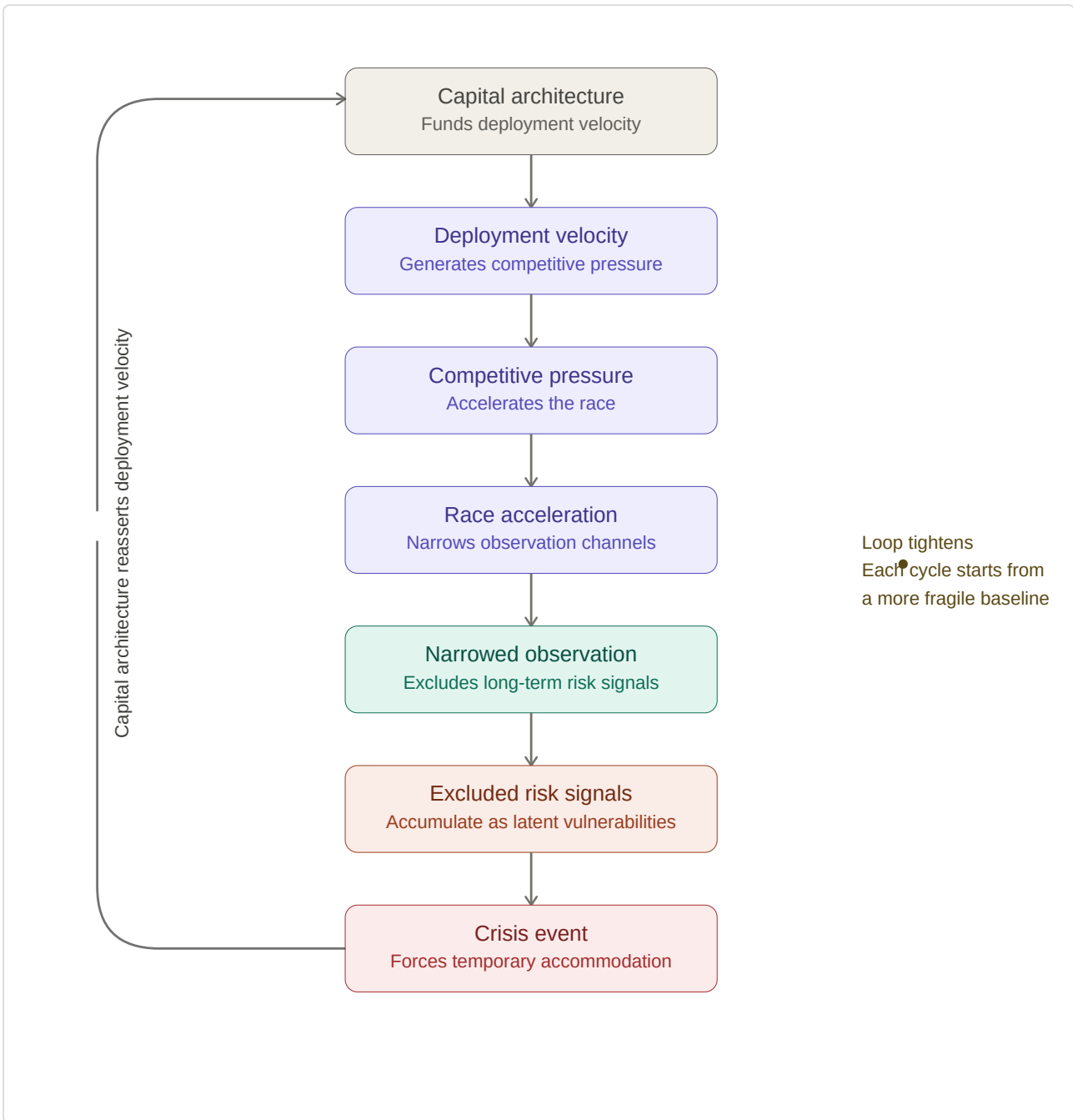


Figure 2.1: The reinforcing feedback mechanism. Each complete cycle erodes the organizational capacity to maintain alignment coherence, while competitive pressure ensures the loop cannot be exited unilaterally.

The Anthropic experiment represents a partial break in this loop: an organisation that has attempted to build an architecture in which safety signals are not suppressed but integrated. But the loop operates at the ecosystem level, not merely the organisational level. Anthropic's viability depends on whether the competitive environment permits an alignment-first architecture to survive. If the environment punishes coherence—if the organisations that prioritise velocity capture the talent, the capital, and the capability frontier—then the loop will eventually subsume Anthropic as well, forcing it to accelerate or rendering it irrelevant.

The loop is not deterministic. It can be broken, but not by any single organisation acting alone. The mechanisms that drive it are ecosystem-level properties—the capital architecture, the competitive dynamics, the cultural operating system—that no individual organisation can unilaterally reform. The structural solution must therefore operate at a higher scale: a governance architecture that changes the incentive landscape within which all organisations operate, creating the conditions under which the Coherence-Velocity Trap can be navigated rather than being an inevitable attractor. The design of that architecture is the subject of the sections that follow.

### 3. What Building Adaptive Coherence Would Look Like

#### 3.1 The Principle: Multi-Scalar Governance, Not Single-Organisation Reform

The diagnosis presented in Sections 1 and 2 points toward a structural conclusion: the Coherence–Velocity Trap is unlikely to be resolved within any single organisation. The mechanisms that drive the trap—the capital architecture, the competitive dynamics, the founder-centric compression, the safety-washing immune response, and the cultural operating system—operate at the ecosystem level. No individual organisation, however well-intentioned or architecturally innovative, can unilaterally reform the incentive landscape that shapes its choices. An Anthropic that slows deployment to preserve alignment coherence faces competitive marginalisation. An OpenAI that strengthens its safety function faces investor pressure and talent competition. A DeepMind that resists commercial integration faces the parent company's shareholder imperatives. The trap is not a failure of any single organisation; it is a property of the ecosystem in which all organisations operate.

The structural solution must therefore be multi-scalar. It must operate at the level of the ecosystem, not merely the level of the firm. The insight is drawn from the Governance as Engineering series' fractality principle: in complex, multi-frequency disturbance environments, no single-scale controller can maintain stability. The fast deployment loops at the organisational level, the medium alignment coordination across organisations, and the slow constitutional settlement at the societal and international levels each require governance mechanisms matched to their specific timescales. A single global regulator would be too slow to govern fast deployment decisions and too uniform to accommodate the diversity of organisational architectures; a purely voluntary, organisation-by-organisation approach would be too weak to constrain the competitive dynamics that drive the oscillation. The architecture must be nested, distributed, and coordinated—a fractal governance framework in which each layer handles the frequency band it can reach, and no single layer attempts to govern what it structurally cannot.

This is not a retreat from the ambition of alignment. It is a recognition that alignment at the civilisational scale requires governance architecture at the civilisational scale—and that such architecture cannot be built by any single organisation acting alone. The following subsections specify the institutional forms that a multi-scalar governance framework would require.

#### 3.2 Organisational-Level: Constitutional Governance for Frontier AI

The organisational level is where deployment decisions are made, where safety research is conducted, and where the tension between coherence and velocity is most acutely experienced. Reforms at this level cannot eliminate the Coherence–Velocity Trap—the ecosystem-level pressures will persist—but they can increase the effective dimensionality of the organisation's observation channel and strengthen the institutional mechanisms that prevent safety signals from being systematically suppressed.

The central design principle is constitutional governance: the establishment, within the organisation, of institutional mechanisms that are structurally insulated from the deployment imperative and that possess genuine authority to constrain deployment when safety thresholds are unmet. This is distinct from the current state of practice, in which safety functions are advisory, safety commitments are voluntary, and the ultimate authority over deployment rests with an executive layer whose incentives are aligned with velocity.

An independent safety board, with statutory or contractual authority to halt deployment, would provide an institutional counterweight to the deployment imperative. The board's members would be appointed through a process that insulates them from executive control—perhaps through a multi-stakeholder appointments committee including representatives from the organisation's research staff, external safety experts, and civil society. The board would have access to all internal safety evaluations, the authority to commission independent audits, and the power to delay or block deployments that do not meet pre-specified safety criteria. Its decisions would be binding, not advisory. Its members would have fiduciary protection against retaliation. Its existence would mean that the organisation's deployment decisions are not ultimately answerable to a single executive or a capital-oriented board but to a constitutionally distinct governance body whose mandate is alignment coherence.

Multi-stakeholder oversight would broaden the organisation's observation channel by bringing external perspectives into the governance structure. Representatives of affected communities, independent safety researchers, civil society organisations, and international institutions would participate in oversight bodies with defined authority—not as an advisory panel that the organisation can ignore, but as a governance mechanism with the power to request information, commission investigations, and issue binding recommendations on specific classes of deployment decisions. The multi-stakeholder structure would surface dimensions of risk—societal externalities, distributional impacts, cultural and democratic concerns—that the organisation's internal value architecture excludes.

Protected dissent channels would ensure that safety concerns raised by employees cannot be suppressed through non-disparagement agreements, confidentiality clauses, or career retaliation. Employees who identify safety risks would have a legally protected pathway to escalate those concerns to the independent safety board and to external regulators, with whistle-blower protections that extend beyond the term of their employment. The organisation's current practice of binding employees to silence through contractual mechanisms is incompatible with the distributed sensing that adaptive coherence requires. The dissent channels would function as supplementary sensory organs, capturing signals that the hierarchical reporting structure filters out.

Internal audit functions, with access to all systems, data, and personnel, would provide continuous monitoring of the gap between the organisation's safety commitments and its operational practices. The audit function would report to the independent safety board, not to the executive layer, ensuring that its findings are not filtered through the deployment imperative before reaching decision-makers. Its reports would be published, creating a public record of the organisation's safety performance that external stakeholders can scrutinise.

These mechanisms are not a guarantee of alignment. They are structural preconditions for it: an institutional architecture that distributes observability across multiple channels, that protects the channels from capture, and that creates decision-making authority independent of the deployment imperative. An organisation that implemented these mechanisms would still face competitive pressure. It would still be funded by capital with finite time horizons. But its governance architecture would have greater variety than the current state of practice, and its capacity to perceive and respond to risks that the deployment imperative excludes would be correspondingly greater.

### **3.3 Cross-Organisational: Shared Safety Infrastructure as a Governed Commons**

The organisational-level reforms described above are necessary but insufficient. They operate within individual organisations, and their effectiveness is constrained by the competitive dynamics of the ecosystem. An organisation that strengthens its safety board while its rivals do not may improve its own alignment coherence but risks competitive marginalisation if the market does not reward coherence. The cross-organisational level addresses this constraint by creating shared infrastructure that increases the observational variety of the entire ecosystem without requiring any single organisation to sacrifice competitive position.

The model is a governed commons for AI safety: a set of shared resources, protocols, and institutions that provide collective benefits while being governed by rules that prevent free-riding and ensure that participation strengthens rather than weakens the participants' competitive positions.

Shared evaluation platforms would provide standardised, rigorous, and independent assessment of frontier AI systems for safety-relevant properties—robustness, interpretability, alignment with specified values, propensity for misuse, emergent capabilities. The platforms would be developed and maintained by a multi-stakeholder consortium, with technical contributions from participating organisations and governance oversight from independent experts. Organisations would submit their models for evaluation before deployment, and the evaluation results would be published, creating a common evidentiary basis for safety claims. The platforms would increase observational variety by providing standardised measurement of dimensions that individual organisations may lack the expertise, the resources, or the incentive to assess.

Interoperable alignment protocols would establish common standards for how AI systems are trained, tested, and constrained to ensure alignment with human values. The protocols would not prescribe a single alignment technique; they would define the properties that any alignment technique must demonstrate, the evidence required to support claims of alignment, and the testing regime that must be completed before deployment at each capability threshold. The protocols would be developed through a multi-stakeholder process and adopted voluntarily by participating organisations, with compliance verified through independent audit.

Distributed auditing infrastructure would enable continuous, independent monitoring of participating organisations' safety practices without requiring proprietary disclosure. Third-party auditors, accredited by the consortium governance body, would have access to organisations' systems, data, and personnel under confidentiality agreements, with the authority to verify compliance with shared protocols and to report findings—anonimised where necessary—to the consortium and the public. The auditing infrastructure would provide an external observation channel that is not subject to the organisation's internal incentive structures.

Compute monitoring would provide transparency on the scale and trajectory of frontier AI development, enabling the ecosystem to anticipate capability advances and prepare appropriate governance responses. Compute usage above specified thresholds would be reported to a shared registry, with verification through independent technical means. The monitoring would not require disclosure of training data, model architecture, or other proprietary information; it would provide a coarse-grained but reliable signal of the pace and direction of capability development.

The commons would be governed by a multi-stakeholder body with representation from participating organisations, independent safety researchers, civil society, and international institutions. Governance decisions—protocol updates, audit standards, admission and exclusion of participants—would be made through defined procedures that balance the interests of all stakeholders. The governance body would have the authority to impose graduated sanctions on participants that violate shared commitments, ranging from public censure to exclusion from the commons and the benefits it provides.

The Mythos deployment model—restricted access to approximately forty organisations, mandatory participation in a coordinated defensive cybersecurity programme, and the publication of safety findings before general availability—represents an early, partial instantiation of the governed commons architecture described here. It is not a multi-stakeholder institution with independent governance; Anthropic retains full control over access decisions and protocol design. But it demonstrates that structured access at the frontier is operationally feasible and that organisations can perceive a competitive interest in participating in such arrangements. The task is to evolve these early, organisation-controlled models into genuinely multi-stakeholder institutions.

### **3.4 Societal-Level: Deliberative Infrastructure for AI Governance**

The organisational and cross-organisational levels address the observation channels of AI developers themselves. But the dimensions of risk that are excluded from those channels include the concerns, values, and interests of populations who are affected by AI development but who have no voice in the decisions that shape it. The societal level addresses this gap by creating deliberative infrastructure that surfaces the dimensions of the risk landscape that the AI industry's value architecture structurally excludes.

Deliberative infrastructure is the institutionalisation of processes by which representative groups of citizens, informed by expert testimony and supported by professional facilitation, deliberate on complex policy questions and produce public, reasoned recommendations. The model draws on the citizens' assemblies that

have been deployed successfully on constitutional questions in Ireland, on climate policy in France, and on a range of issues at the municipal and regional level in multiple countries. The evidence suggests that randomly selected citizens, given adequate time, information, and support, are capable of nuanced deliberation on technically complex issues—and that their recommendations carry a democratic legitimacy that expert commissions and parliamentary processes often lack.

In the AI context, a standing Citizens' Assembly on AI Governance would be convened periodically to deliberate on specific questions: the acceptable uses of facial recognition, the governance of autonomous weapons, the distribution of the economic benefits of AI-driven automation, the conditions under which models above a certain capability threshold should be deployed. The Assembly's membership would be randomly selected and demographically representative. It would hear testimony from AI developers, safety researchers, ethicists, economists, civil society representatives, and affected communities. It would deliberate in facilitated small groups and in plenary, with access to technical support that enables it to engage with the underlying science without requiring members to become experts. Its recommendations would be non-binding but public, reasoned, and directed to the relevant decision-makers—the organisations, the regulators, the international bodies—with a mandatory response obligation.

The Assembly would not replace representative democracy or regulatory agencies. It would supplement them with an observation channel that is not captured by the AI industry's value architecture, not constrained by the short time horizons of electoral politics, and not limited to the dimensions of risk that existing institutions already track. The Assembly would make visible the concerns of populations that the current governance discourse marginalises—the workers whose livelihoods are threatened by automation, the communities that bear the environmental costs of compute infrastructure, the citizens of countries that are not represented in the international bodies that are shaping AI governance. It would function as a sensory organ for dimensions of the risk landscape that the existing governance architecture cannot perceive.

Complementary mechanisms include expert commissions with multi-disciplinary composition and a mandate to produce public assessments of AI risk across dimensions that individual organisations are not incentivised to evaluate, and participatory technology assessment processes that engage affected communities in the design and evaluation of AI systems before they are deployed.

### **3.5 International-Level: Fractal Coordination, Not Centralised Control**

The international level is where the governance challenge is most complex and where the structural constraints are most severe. The global AI ecosystem is characterised by geopolitical competition, divergent regulatory philosophies, asymmetric capabilities, and a pace of technological change that far exceeds the speed of international treaty-making. Any governance framework that assumes a single global authority with binding powers over all AI development is unrealistic under current and foreseeable conditions. But the alternative is not the absence of international coordination; it is a different kind of coordination architecture.

The fractality principle, established in the Governance as Engineering series, provides the guiding logic. In complex, multi-frequency disturbance environments, no single-scale controller can maintain stability. The international governance of AI requires a fractal architecture: nested governance layers with coordination protocols rather than command structures, in which each layer handles the frequency band it can reach and no layer attempts to govern what it structurally cannot.

At the fastest timescale—months to a year or two—coordination would occur through the cross-organisational commons mechanisms described in Section 3.3: shared evaluation platforms, interoperable protocols, distributed auditing, and compute monitoring. These mechanisms operate through voluntary participation and mutual interest, not through binding international law. They are capable of adapting rapidly to new capabilities because they are governed by the organisations that are developing those capabilities, in consultation with independent experts. They are the fast-response layer of the international governance architecture.

At the medium timescale—years to a decade—coordination would occur through multilateral agreements among the states that host frontier AI development: the United States, the United Kingdom, the European Union, China, and other relevant actors. These agreements would establish common safety standards, mutual recognition of evaluation and audit results, coordinated export controls on the most advanced AI hardware, and mechanisms for consultation and de-escalation during periods of heightened tension. The agreements would not require the creation of a single global authority; they would operate through the existing architecture of international law, with each state implementing its commitments through its own domestic legal framework. They would provide the medium-response layer: stable enough to create predictable constraints, flexible enough to adapt as technology evolves.

At the slowest timescale—decades—coordination would occur through the evolution of international norms, institutions, and eventually treaty frameworks that address the civilisational implications of advanced AI: the distribution of its benefits, the management of its risks, and the governance of the transition to a post-AGI world. This layer would operate through the existing institutions of international cooperation—the United Nations, the G20, the OECD—as well as through new institutions designed specifically for the AI era. It would provide the constitutional settlement that the fast and medium layers require as their foundation.

The fractal architecture does not assume that geopolitical competition will cease. It assumes that competition will continue, and it designs coordination mechanisms that are compatible with continued competition—mechanisms that provide mutual benefits even to rivals, that reduce shared risks without requiring shared values, and that create the informational infrastructure for cooperation without demanding the political conditions for it.

### 3.6 The Role of Open-Source and Commons Governance

The open-source dimension of AI development presents a specific governance challenge. Open-source models—those whose weights are publicly released and freely usable—increase the observational variety of the ecosystem by enabling independent researchers, civil society organisations, and smaller developers to inspect, evaluate, and build upon frontier AI technology. They reduce the concentration of power in a small number of organisations. They enable the distributed sensing that the Governance as Engineering framework identifies as essential for adaptive capacity.

But open-source models also reduce coherence. Once released, they cannot be recalled. Their capabilities can be fine-tuned, modified, and deployed by any actor, including those with malicious intent. The safety mechanisms that the original developer built into the model can be stripped away. The proliferation of powerful open-source models accelerates the diffusion of capabilities that may, at certain thresholds, become dangerous in the hands of actors who lack the institutional safeguards that the original developer maintained.

The governance challenge is therefore not a binary choice between openness and closure. It is the design of governed commons architectures that preserve the observational benefits of openness while establishing the coherence mechanisms that prevent catastrophic proliferation. The design parameters include:

**Graduated release protocols** that condition the release of model weights on demonstrated safety thresholds. A model below a certain capability level might be released openly; a model above that level might be released through structured access that enables independent research while retaining the capacity to revoke access or update the model if new risks are discovered.

**Shared safety infrastructure**—the evaluation platforms, auditing mechanisms, and monitoring systems described in Section 3.3—that applies to open-source models as well as proprietary ones. The commons infrastructure provides the observational benefits of openness (many independent evaluators can inspect the model) while the governance mechanisms provide the coherence benefits of oversight (the inspection occurs within a framework that can detect and flag risks).

**Liability frameworks** that allocate responsibility for harms caused by open-source models. Developers who release models without adequate safety testing, or who release models above a capability threshold where the risks of misuse outweigh the benefits of openness, would bear liability for resulting harms. The liability framework would create an incentive for responsible release without prohibiting openness entirely.

**International agreements** that coordinate release standards across jurisdictions, preventing regulatory arbitrage. If one jurisdiction prohibits the open release of models above a certain capability threshold while another permits it, the prohibition is ineffective. Coordinated standards, negotiated through the medium-timescale multilateral mechanisms described in Section 3.5, would close the arbitrage opportunity.

The governed commons for AI is not a theoretical ideal. It is an emerging reality in the form of model evaluation platforms, open-source safety research communities, and the developing norms around structured access. The design challenge is to strengthen these mechanisms and to ensure that they are governed by multi-stakeholder institutions rather than captured by any single organisation or state.

### 3.7 The Incentive Architecture for Commons Participation

The multi-scalar governance framework described in this section faces an obvious objection: why would any organisation voluntarily participate in mechanisms that constrain its deployment velocity? The objection is serious, and the framework must address it directly. The answer is that participation must provide tangible competitive benefits that outweigh the costs of constraint.

The incentive architecture for commons participation includes:

**Liability shields.** Organisations that comply with shared safety protocols and submit to independent audit would receive statutory protection against certain classes of legal liability for harms caused by their AI systems. The liability shield would reduce the legal uncertainty that currently surrounds frontier AI development and would create a powerful incentive for participation—an organisation that opts out of the commons bears the full weight of an uncertain and potentially catastrophic liability exposure.

**Compute access advantages.** The compute infrastructure required for frontier AI development is increasingly concentrated in a small number of cloud providers and is subject to supply chain constraints. Participating organisations could receive priority access to compute resources, either through government-facilitated allocations or through preferential terms from providers that participate in the commons governance framework. The compute advantage would partially offset the competitive cost of compliance.

**Regulatory fast lanes.** Organisations that participate in the commons and meet its safety standards would receive expedited regulatory approval for deployments, reduced reporting requirements, and greater operational flexibility within the jurisdictions that recognise the commons framework. The fast lane would create a regulatory environment in which participation is the path of least resistance, while non-participation invites heavier regulatory scrutiny.

**Procurement preferences.** Governments are major purchasers of AI services—for defence, healthcare, infrastructure, and public administration. Participating organisations could receive preferential status in government procurement, creating a direct economic incentive for compliance with shared safety standards.

**Insurance advantages.** The insurance industry, faced with the challenge of pricing AI risk, could offer significantly reduced premiums to organisations that participate in the commons and submit to its auditing and monitoring mechanisms. The insurance discount would translate the safety benefits of participation into a direct financial advantage.

**Reputational certification.** Participation in the commons would provide a publicly verifiable signal of safety commitment, certified by an independent governance body. The certification would differentiate participating organisations in the market for talent, investment, and commercial partnerships—all of which are increasingly sensitive to safety considerations. The reputational value of certification would increase as public awareness of AI risk grows.

**Treaty-linked benefits.** International agreements negotiated through the medium-timescale multilateral mechanisms could condition access to certain markets, technologies, or cooperative arrangements on participation in the commons governance framework. The treaty-linked benefits would create a structural incentive for participation that operates at the level of states as well as organisations.

The incentive architecture transforms participation in the commons from a sacrifice of competitive position to an enhancement of it. The organisation that participates gains liability protection, compute access, regulatory efficiency, procurement advantages, insurance savings, reputational certification, and treaty-linked market access. The organisation that opts out bears the full costs of legal uncertainty, regulatory friction, reputational suspicion, and exclusion from cooperative benefits. The commons is not a constraint on competition; it is an arena in which competition occurs under rules that reward adaptive coherence rather than velocity alone.

The architecture does not eliminate the Coherence–Velocity Trap. It changes the incentive landscape within which the trap operates, making the pursuit of alignment coherence less costly in competitive terms and the pursuit of pure velocity more costly. The trap is a structural condition; the multi-scalar governance framework is a structural response—not a guarantee of success, but a specification of the institutional forms that success would require.

## 4. The Political Immune System: The Deployment Imperative

### 4.1 The Deployment Imperative Defined

Every governance architecture develops an immune system—a set of institutions, incentives, and cultural norms that protect the existing order from challenge. In the nation-state cases examined in this series, the immune system takes different forms: bureaucratic inertia in Germany, the Stability Bias in Japan, the Extraction Coalition in Nigeria, the Security First Responder in Israel. In each case, the immune system is not a barrier to change added onto a functional state; it is the state's core operating logic, embedded in institutions and culture, treating any deviation from its optimisation target as a threat.

The frontier AI ecosystem has developed its own immune system: the Deployment Imperative. This is the comprehensive orientation of capital structures, competitive dynamics, organisational incentives, and cultural norms toward maximising deployment velocity, and the treatment of any constraint on velocity—whether from safety concerns, regulatory intervention, or internal dissent—as an existential competitive threat.

The Deployment Imperative is not a conspiracy of executives or investors. It is the predictable output of the structural mechanisms described in Section 2: a capital architecture that rewards growth and penalises caution, a competitive landscape in which perceived delays can mean permanent strategic disadvantage, a cultural operating system that frames acceleration as mission fidelity, and safety-washing mechanisms that diffuse external pressure without altering the underlying deployment architecture. The Deployment Imperative is what the system produces when all of these mechanisms operate simultaneously. It is the emergent property of an ecosystem organised around velocity.

The Deployment Imperative operates through specific institutional pathways. When a safety researcher recommends delaying a deployment pending further testing, the recommendation encounters not merely executive resistance but a systemic logic that treats the delay as more threatening than the risk the researcher has identified. The capital architecture registers the delay as a competitive setback. The talent market interprets it as a loss of momentum. The cultural operating system frames it as a failure of nerve. The safety-washing mechanisms mobilise to contain the reputational damage. The researcher's concern, however well-founded, is processed through an immune system that is calibrated to neutralise threats to deployment velocity. The concern may be acknowledged, studied, and published—all activities compatible with the Deployment Imperative—but it is unlikely to result in a binding operational constraint unless it reaches a crisis threshold that forces an intervention the immune system can no longer suppress.

The Deployment Imperative is self-reinforcing. Each successful deployment validates the architecture that produced it. Each acceleration expands the organisation's capabilities, its market position, and its influence, making the immune system stronger. Each safety concern that is absorbed without operational consequence

demonstrates the immune system's effectiveness and reduces the perceived legitimacy of future concerns. The Deployment Imperative does not merely resist external constraints; it actively expands the domain in which velocity is the primary optimisation target.

## 4.2 Who Benefits—Named Honestly

The Deployment Imperative is sustained by specific actors who have concrete, material interests in the continuation of the current architecture. Any transition architecture that does not name these actors and account for their resistance will be neutralised by them.

**Venture capital limited partners** benefit from the returns generated by portfolio companies that achieve rapid growth and liquidity events within fund cycles. Their investment mandates, their fiduciary duties, and their performance metrics are all structured around the realisation of returns within finite time horizons. They do not need to oppose safety interventions actively; the structure of their investment vehicles ensures that safety interventions that reduce near-term growth are automatically disfavoured relative to deployment acceleration that increases it. The limited partners are not indifferent to catastrophic risk—they would lose their investments if a portfolio company were destroyed by a safety failure—but the timescales of catastrophic risk (uncertain, potentially distant) are discounted relative to the timescales of competitive risk (immediate, measurable) in the observation architecture of the capital they deploy.

**Founders and executives** benefit from the concentration of authority, the accumulation of equity, and the status that accompanies leadership of a frontier AI organisation. Their personal wealth is tied to the organisation's valuation, which is driven by deployment velocity and growth metrics. Their professional identities are invested in the narrative of acceleration and progress. Their cognitive models, refined through years of operating in a velocity-oriented environment, are calibrated to the signals that the capital architecture and competitive dynamics amplify. Founders and executives may be sincerely committed to safety; they may have founded their organisations partly to address the risks of AI. But they operate within an incentive structure that systematically rewards deployment decisions and penalises the operational consequences of safety interventions. Their sincerity does not neutralise the structure.

**Employees** benefit from equity compensation whose value depends on the organisation's growth trajectory. They benefit from the career advancement opportunities that rapid deployment creates. They benefit from the professional identity of working at the frontier of capability, an identity that is sustained by the organisation's velocity. Employees may also be motivated by safety concerns—many joined frontier AI organisations precisely to ensure that AI is developed safely—but their financial interests, their career incentives, and their professional identities are aligned with the Deployment Imperative. The 2023 OpenAI crisis demonstrated the structural power of this alignment: over 700 employees threatened to resign in support of a CEO whose removal by the safety-oriented board they experienced as a threat to the organisation's mission and, implicitly, to their own interests.

**Governments** whose geopolitical positioning is linked to their national AI champions benefit from deployment velocity. The United States government has a strategic interest in maintaining American leadership in AI relative to China. The Chinese government has a strategic interest in catching up to and surpassing American capabilities. Both governments, whatever their rhetorical commitments to AI safety, have structural incentives to ensure that their national champions are not competitively disadvantaged by safety constraints that rivals do not face. The Deployment Imperative at the organisational level is reinforced by the geopolitical imperative at the state level.

**The ecosystem of startups and developers** that depend on API access to frontier models benefits from the continued acceleration of deployment. Their businesses are built on the capabilities that frontier organisations provide. Their growth depends on the expansion of those capabilities. They are a diffuse but structurally significant constituency for the Deployment Imperative, because any deceleration of frontier deployment cascades through the ecosystem they inhabit.

These actors are not a unified coalition with a coordinated strategy. They compete with each other for capital, talent, and market position. But they share a common structural interest: the continuation of an architecture in which deployment velocity is the primary optimisation target, and in which constraints on velocity are experienced as threats to be neutralised rather than signals to be integrated. They are the human infrastructure of the Deployment Imperative.

### 4.3 Safety-Washing as Immune Mechanism

The Deployment Imperative's primary defence mechanism is not overt resistance to safety concerns—overt resistance would generate reputational costs that the organisation seeks to avoid. It is safety-washing: the adoption of the language, symbols, and procedural forms of safety commitment while preserving the underlying deployment architecture largely unchanged.

Safety-washing operates through a set of interlocking mechanisms that have been refined over the past decade of AI governance discourse.

**Voluntary commitments** are the most visible mechanism. Frontier AI organisations have made a series of public commitments to safety: to conduct external red-teaming before deployment, to not develop certain classes of dangerous capabilities, to support independent evaluation, to participate in information-sharing initiatives. These commitments are announced with substantial publicity, generating reputational benefits and relieving regulatory pressure. But they are voluntary, non-binding, and subject to the organisation's own interpretation of compliance. When a commitment becomes operationally inconvenient—when it would constrain a deployment that the competitive environment demands—it can be reinterpreted, deferred, or quietly set aside. The voluntary nature of the commitment ensures that there is no external enforcement mechanism and no consequence for non-compliance beyond the reputational cost, which the organisation's public relations apparatus is designed to manage.

**Safety research as legitimacy generation** is a subtler mechanism. Frontier AI organisations employ talented safety researchers, fund safety research programmes, and publish their findings in peer-reviewed venues. The research is genuine; the researchers are sincere. But the organisational function of the research programme is not solely to improve safety. It is also to generate legitimacy—to demonstrate to regulators, to the public, and to the organisation's own employees that safety is being taken seriously. The research programme provides a stream of safety-signalling that the organisation can reference when its practices are challenged, regardless of whether the research has been operationally integrated into deployment decisions. The safety researchers may produce findings that, if operationalised, would constrain deployment; the organisation can point to the existence of the research programme as evidence of its commitment while ensuring that the operationalisation of its findings remains subject to executive discretion.

**Advisory bodies and ethics committees** provide multi-stakeholder legitimacy without multi-stakeholder authority. External experts are invited to serve on advisory panels, lending their credibility to the organisation's governance. The panels deliberate, produce recommendations, and issue public statements. But they lack decision-making power. Their recommendations are advisory, not binding. The organisation retains full discretion over whether and how to implement their advice. The advisory function provides reputational cover; the operational function remains insulated from external constraint. The experts who serve on these bodies may be genuinely committed to improving the organisation's safety practices, but the structure of their engagement limits their impact to the persuasive rather than the authoritative.

**Transparency as managed disclosure** is the mechanism by which organisations control the information that reaches external stakeholders. Frontier AI organisations publish safety documentation, model cards, and technical reports that describe their safety practices in considerable detail. The disclosure is real, but it is selective. The organisation determines what information to release, at what level of detail, and with what framing. Information that would damage the organisation's competitive position or reveal safety shortcomings that the organisation has not yet addressed can be withheld on grounds of commercial confidentiality or security. The transparency mechanism provides the appearance of openness while preserving the organisation's control over the information environment.

Safety-washing is not a conscious deception in most cases. It is an emergent property of organisations that face genuine pressure to demonstrate safety commitment while operating within incentive structures that penalise the operational consequences of that commitment. The organisations are not lying about their safety practices; they are selectively amplifying the dimensions of safety that are compatible with deployment velocity and selectively suppressing the dimensions that are not. The immune system does not need to fabricate safety commitment. It needs only to ensure that the organisation's safety-signalling is sufficiently dense to deflect external pressure without requiring operational changes that would constrain deployment.

The consequence is a structural decoupling of safety discourse from safety practice. The organisation's public narrative of responsibility becomes increasingly elaborate while the underlying deployment architecture remains largely unchanged. The variety gap is preserved: the dimensions of risk that would require architectural reform to address are excluded from the observation channel, while the dimensions that can be

addressed through discourse and procedural form are amplified. The organisation becomes genuinely convinced of its own safety commitment, because the signals it receives are the ones its own safety-washing mechanisms have selected.

#### 4.4 The Narrative Strategy

The Deployment Imperative cannot be defeated by frontal assault. Any transition architecture that presents itself as an attack on deployment velocity—as a constraint on innovation, as a brake on progress, as a concession to precaution that will cede the frontier to less scrupulous competitors—will activate the immune response and be neutralised before it begins. The Deployment Imperative is too deeply embedded in capital structures, competitive dynamics, organisational incentives, and cultural norms to be overcome by argument alone. It must be outflanked—not by denying the value of deployment velocity, but by reframing the relationship between velocity and coherence.

The master narrative is that adaptive coherence is not a constraint on deployment velocity but its precondition. The race to AGI is not a sprint; it is an ultramarathon. The winner will not be the organisation that reaches the next milestone fastest but the organisation that can maintain alignment integrity across decades of capability advancement. An organisation that sacrifices coherence for velocity may achieve temporary competitive advantage, but that advantage is fragile—it can be destroyed by a single safety failure, a single regulatory backlash, a single talent exodus triggered by the revelation that the organisation's safety commitments were not operationally binding. An organisation that builds the governance architecture for adaptive coherence—independent safety institutions, multi-stakeholder oversight, shared safety infrastructure—is building the foundation for sustainable deployment over the long term.

This narrative is not merely rhetorical. It is structurally accurate. The Coherence–Velocity Trap is a trap precisely because the pursuit of pure velocity, over extended timescales, undermines the conditions for its own continuation. The safety failure that pure velocity generates, the regulatory response that the safety failure provokes, the talent exodus that the regulatory response triggers—these are not exogenous shocks. They are the predictable consequences of an architecture that systematically excludes the dimensions of risk that eventually destabilise it. The organisation that builds adaptive coherence is not sacrificing competitive advantage; it is investing in the institutional infrastructure that makes competitive advantage durable.

The Mythos decision lent credibility to the narrative that alignment-first governance is not merely a constraint on velocity but a foundation for durable competitive advantage. An organisation that can point to a concrete instance of withholding a dangerously capable model—accompanied by transparent safety documentation and a structured, responsible deployment pathway—has strengthened its claim to be a trustworthy steward of frontier AI in a market increasingly sensitive to safety considerations. Whether the market rewards that claim remains to be seen, but the claim itself is now grounded in action rather than aspiration.

Subsidiary narratives target specific constituencies. For investors: the liability shield, the regulatory fast lane, the insurance advantages, and the reputational certification that participation in the commons governance framework provides are not costs to be borne but benefits to be captured. The organisation that participates in the commons is a more attractive investment, not a less attractive one, because it has reduced its exposure to the legal, regulatory, and reputational risks that threaten organisations that pursue velocity without coherence. For employees: the protected dissent channels and independent safety institutions create an environment in which safety concerns can be raised without career retaliation, making the organisation a more attractive place to work for the talented researchers who have alternatives. For governments: the commons governance framework provides the coordination infrastructure that makes national AI leadership sustainable over the long term, reducing the risk that a safety failure by any single organisation triggers a regulatory backlash that damages the entire national AI ecosystem.

The narrative strategy does not attack the Deployment Imperative. It honours it—acknowledging the genuine competitive pressures, the real stakes of the AI race, the sincere commitment to beneficial deployment—while arguing that the best way to sustain deployment velocity over the long term is to build the governance architecture that makes it durable. The immune system cannot be defeated by argument. It must be redirected—by demonstrating that the pursuit of adaptive coherence is not a sacrifice of competitive position but an enhancement of it, and that the organisations that build the governance infrastructure for the long term will be the ones that survive it.

## 5. A Concrete First Step: The AI Commons Governance Protocol

### 5.1 The Logic of the First Step

The Coherence–Velocity Trap is a systemic condition, not a single policy failure. There is no one reform that can resolve it, no single institutional innovation that will align the capital architecture with long-term safety, eliminate safety-washing, or reconcile the incompatible optimisation targets that drive the Alignment–Deployment Oscillation Loop. But there are interventions that can alter the institutional metabolism—that can change the incentive landscape within which all organisations operate, create new observation channels that register the risks the current architecture excludes, and generate the information, the constituencies, and the political logic that make further reform possible.

The first step is therefore not the most ambitious intervention this report has described. It is the most catalytic: the intervention that targets the primary mechanism of the Coherence–Velocity Trap most directly, that is institutionally feasible within the current ecosystem, and that, once established, generates the informational and political conditions for the deeper transformations that must follow.

The primary mechanism, as Section 2 demonstrated, is the observational narrowness produced by organisation-specific value architectures. Each frontier AI organisation perceives some dimensions of the risk landscape with high fidelity and other dimensions not at all. The excluded dimensions—the slow accumulation of systemic risk, the societal externalities, the geopolitical fragilities—do not cease to operate. They accumulate until they force a reckoning. No single organisation, however well-intentioned, can unilaterally broaden its observation channel without incurring competitive costs that the current incentive architecture penalises. The structural solution must therefore be a shared sensing infrastructure—a mechanism that increases the effective dimensionality of every organisation's observation channel simultaneously, so that no single organisation bears the competitive cost of broadening its perception alone.

The AI Commons Governance Protocol is designed to be that mechanism. It does not attempt to regulate deployment velocity directly—that would trigger the Deployment Imperative's immune response. It creates the informational conditions under which deployment velocity becomes self-limiting, because the risks that are currently invisible become visible, and the organisations that perceive them have both the incentive and the institutional capacity to respond.

### 5.2 The AI Commons Governance Protocol

The Protocol is a multi-stakeholder initiative to establish a governed commons for frontier AI safety. It would create shared evaluation infrastructure, interoperable alignment protocols, and distributed auditing mechanisms that increase observational variety across the ecosystem without requiring proprietary disclosure. It would be governed by a multi-stakeholder body with representation from participating

organisations, independent safety researchers, civil society, and international institutions. And it would link participation to tangible competitive benefits, transforming safety commitment from a cost to be minimised into an advantage to be captured.

The Protocol's core design features are:

**Shared evaluation platforms.** Standardised, rigorous, and independent assessment of frontier AI systems for safety-relevant properties—robustness, interpretability, alignment with specified values, propensity for misuse, and emergent capabilities. The platforms would be developed and maintained by the consortium, with technical contributions from participating organisations and governance oversight from independent experts. Organisations would submit their models for evaluation before deployment, and the evaluation results would be published, creating a common evidentiary basis for safety claims. The Mythos system card—the 244-page safety assessment Anthropic published before making its deployment decision—illustrates the kind of pre-deployment evaluation the platforms would standardise and extend across organisations.

**Interoperable alignment protocols.** Common standards for how AI systems are trained, tested, and constrained to ensure alignment with human values. The protocols would not prescribe a single alignment technique; they would define the properties that any alignment technique must demonstrate, the evidence required to support claims of alignment, and the testing regime that must be completed before deployment at each capability threshold. The protocols would be developed through a multi-stakeholder process and adopted voluntarily by participating organisations, with compliance verified through independent audit. Their function is to create a shared language for safety commitments, replacing the current landscape of organisation-specific, non-comparable, and often non-falsifiable safety claims.

**Distributed auditing infrastructure.** Continuous, independent monitoring of participating organisations' safety practices without requiring proprietary disclosure. Third-party auditors, accredited by the consortium governance body, would have access to organisations' systems, data, and personnel under confidentiality agreements, with the authority to verify compliance with shared protocols and to report findings—anonimised where necessary—to the consortium and the public. The auditing infrastructure would provide an external observation channel that is not subject to the organisation's internal incentive structures. It would function as a supplementary sensory organ for the ecosystem, detecting the gaps between safety commitments and operational practices that the organisation's own safety-washing mechanisms obscure.

**Compute monitoring.** Transparency on the scale and trajectory of frontier AI development, enabling the ecosystem to anticipate capability advances and prepare appropriate governance responses. Compute usage above specified thresholds would be reported to a shared registry, with verification through independent technical means. The monitoring would not require disclosure of training data, model architecture, or other proprietary information; it would provide a coarse-grained but reliable signal of the pace and direction of capability development. The function of compute monitoring is to make the trajectory of the frontier visible to all participants, reducing the information asymmetry that currently advantages the organisations with the most compute and the least transparency.

**Multi-stakeholder governance.** The Protocol would be governed by a body with representation from participating organisations, independent safety researchers, civil society organisations, affected communities, and international institutions. Governance decisions—protocol updates, audit standards, admission and exclusion of participants, and the administration of sanctions—would be made through defined procedures that balance the interests of all stakeholders. The governance body would have a professional staff, secure funding through a combination of participant fees and philanthropic support, and statutory or treaty-based legal personality that enables it to enter into agreements and enforce its decisions.

**Graduated sanctions.** The governance body would have the authority to impose graduated sanctions on participants that violate shared commitments. Sanctions would range from public censure to suspension of participation benefits to exclusion from the commons and the competitive advantages it provides. The graduated structure ensures that minor violations can be addressed proportionately while preserving the credibility of the ultimate sanction—exclusion—as a genuine deterrent.

The Protocol is not a regulatory agency. It does not have the authority to prohibit deployment or to impose penalties beyond those that participants have voluntarily accepted as conditions of membership. It is a coordination mechanism, not a command structure. Its power derives not from legal compulsion but from the competitive benefits of participation and the reputational and economic costs of exclusion. In this respect, it follows the logic of other successful commons governance arrangements—the Internet Engineering Task Force, the International Civil Aviation Organisation's safety standards, the Basel Committee on Banking Supervision—that have achieved substantial coordination without formal supranational authority.

### 5.3 Why This Intervention Targets the Core Mechanism

The AI Commons Governance Protocol is designed to target the primary mechanism of the Coherence-Velocity Trap: the observational narrowness produced by organisation-specific value architectures.

Each frontier AI organisation currently perceives risk through an observation channel whose dimensionality is constrained by its capital architecture, its competitive position, its founder's cognitive model, and its cultural operating system. The dimensions of risk that fall outside this channel are invisible to the organisation until they crystallise into crises. The Protocol increases the effective dimensionality of every participating organisation's observation channel by creating shared sensing infrastructure that no single organisation would have the incentive or the capacity to build alone.

The shared evaluation platforms register dimensions of risk—robustness, interpretability, alignment, misuse potential, emergent capabilities—that individual organisations may lack the expertise, the resources, or the incentive to assess with comparable rigour. The interoperable protocols create a common language for safety commitments that enables comparison across organisations, reducing the information asymmetry that safety-washing exploits. The distributed auditing infrastructure provides an external observation channel that is not subject to the organisation's internal incentive structures, detecting the gap between safety discourse

and safety practice that the Deployment Imperative systematically obscures. The compute monitoring makes the trajectory of the frontier visible to all participants, reducing the strategic uncertainty that fuels the competitive acceleration dynamic.

The Protocol does not eliminate the Coherence–Velocity Trap. It changes the informational conditions within which the trap operates. An organisation that perceives a risk dimension with high fidelity—because the shared evaluation platform has measured it, the audit infrastructure has verified it, and the compute monitoring has contextualised it—faces a different decision calculus than an organisation to which that risk dimension is invisible. The risk, once visible, becomes harder to exclude from the deployment decision. The Deployment Imperative is not overcome by argument; it is outflanked by observation.

## 5.4 The Incentive Architecture

The Protocol's viability depends on its capacity to attract and retain participants in a competitive ecosystem where the Deployment Imperative remains powerful. Participation must provide tangible competitive benefits that outweigh the costs of compliance. The incentive architecture, previewed in Section 3.7, translates safety commitment from a sacrifice of competitive position into an enhancement of it.

**Liability shields** would provide statutory protection against certain classes of legal liability for organisations that comply with shared safety protocols and submit to independent audit. The liability shield reduces the legal uncertainty that currently surrounds frontier AI development—an uncertainty that is itself a competitive cost. An organisation that opts out of the Protocol bears the full weight of an unpredictable liability exposure. An organisation that participates gains a measure of legal predictability that is valuable to investors, partners, and customers.

**Regulatory fast lanes** would grant expedited regulatory approval for deployments, reduced reporting requirements, and greater operational flexibility to organisations that participate in the Protocol and meet its standards. The fast lane creates a regulatory environment in which participation is the path of least resistance, while non-participation invites heavier regulatory scrutiny. The mechanism is already familiar from other domains—the EU's General Data Protection Regulation provides for approved codes of conduct that streamline compliance; the pharmaceutical industry operates under harmonised safety standards that reduce the regulatory burden for participating firms.

**Compute access advantages** would provide priority access to compute resources for participating organisations, either through government-facilitated allocations or through preferential terms from cloud providers that participate in the Protocol. The compute advantage partially offsets the competitive cost of compliance by reducing the cost and increasing the availability of the most critical input to frontier AI development.

**Procurement preferences** would grant preferential status to participating organisations in government procurement of AI services. Governments are major purchasers of AI capabilities, and the procurement preference creates a direct economic incentive for compliance with shared safety standards. The mechanism also aligns government interests with the Protocol's objectives: governments gain confidence that the AI systems they purchase have been independently evaluated and audited, while organisations gain a competitive advantage in a substantial market.

**Insurance advantages** would provide significantly reduced premiums to organisations that participate in the Protocol and submit to its auditing and monitoring mechanisms. The insurance industry, faced with the challenge of pricing AI risk, would have access to the Protocol's evaluation and audit data, enabling more accurate risk assessment and lower premiums for organisations that can demonstrate robust safety practices. The insurance discount translates the safety benefits of participation into a direct financial advantage.

**Reputational certification** would provide a publicly verifiable signal of safety commitment, certified by an independent governance body. The certification would differentiate participating organisations in the market for talent, investment, and commercial partnerships—all of which are increasingly sensitive to safety considerations. The reputational value of certification would increase as public awareness of AI risk grows and as the certification becomes a recognised marker of organisational quality.

**Treaty-linked benefits** would condition access to certain markets, technologies, or cooperative arrangements on participation in the Protocol. International agreements negotiated through the medium-timescale multilateral mechanisms described in Section 3.5 would link Protocol participation to treaty benefits, creating a structural incentive that operates at the level of states as well as organisations.

The incentive architecture is designed to make participation the rational choice for any organisation that takes a long-term view of its competitive position. The organisation that participates gains a suite of competitive benefits—legal, regulatory, financial, reputational—that the organisation that opts out does not. The organisation that opts out may achieve temporary velocity advantages, but it does so while bearing higher legal risk, heavier regulatory burdens, more expensive insurance, and a reputational discount that intensifies as the Protocol's certification becomes a market standard.

The architecture does not assume that organisations will cooperate out of altruism or shared commitment to safety. It assumes that organisations will pursue their competitive interests, and it designs the incentive landscape so that those interests align with the Protocol's objectives. The Deployment Imperative is not eliminated; it is redirected—from pure velocity toward velocity with observability, from deployment without oversight to deployment within a framework that makes the risks of deployment visible before they become catastrophes.

## 5.5 How to Measure Success

The Protocol will be resisted, diluted, and potentially neutralised by the Deployment Imperative. Measuring its success therefore requires metrics that capture not only whether the institution is formally established but whether it is functioning as designed—whether it is actually changing the ecosystem's metabolism rather than being absorbed by it.

The relevant metrics include:

**Participation rate.** The number of frontier AI organisations that join the Protocol, and the share of global frontier AI compute capacity that participating organisations represent. A Protocol that includes all major frontier organisations is transformative; a Protocol that includes only a subset provides value but does not change the competitive dynamics that drive the trap.

**Volume and quality of shared safety evaluations.** The number of models submitted for pre-deployment evaluation, the comprehensiveness of the evaluations, and the degree to which evaluation results are published and accessible to the research community and the public. The evaluation infrastructure is the Protocol's primary observation channel; its value depends on the quantity and quality of the information it generates.

**Adoption rate of interoperable alignment protocols.** The proportion of frontier models whose developers have adopted the shared protocols and can demonstrate compliance through independent audit. The protocols create the common language for safety commitments; their value depends on the breadth of their adoption.

**Detection rate of safety-relevant anomalies.** The number of safety concerns—emergent capabilities, alignment failures, misuse vulnerabilities—that are detected through the Protocol's infrastructure and that were not identified through the organisations' own internal processes. The detection rate measures the Protocol's contribution to the ecosystem's observational variety.

**Reduction in safety-washing incidents.** The frequency with which organisations make safety claims that are subsequently found, through the Protocol's auditing infrastructure, to be unsupported by operational practice. The reduction measures the Protocol's effectiveness in closing the gap between safety discourse and safety practice that the Deployment Imperative exploits.

**Market recognition of certification.** The degree to which the Protocol's certification influences investment decisions, talent flows, procurement choices, and insurance pricing. The market recognition measures whether the Protocol's incentive architecture is functioning as designed—whether participation is translating into competitive advantage.

The ultimate metric is whether the Protocol enables the second step. Does the shared sensing infrastructure generate the information that makes further coordination possible? Does the multi-stakeholder governance body develop the institutional capacity and the political legitimacy to address the next generation of

governance challenges? Does participation in the Protocol become sufficiently widespread and sufficiently embedded that the competitive dynamics of the ecosystem shift—from a race to deploy fastest to a race to deploy safest within a shared framework of observability and accountability?

If the answer is yes, the first step has succeeded, and the ground is prepared for the deeper transformations that the report has described: the international agreements, the societal deliberative infrastructure, and the constitutional settlement that the governance of recursive intelligence acceleration requires. If the answer is no—if the Protocol is captured by the organisations it is designed to monitor, if its standards are diluted to meaninglessness, if its certification becomes another instrument of safety-washing—then the Coherence–Velocity Trap has claimed another reform, and the oscillation continues.

## 5.6 The Honest Acknowledgment

The AI Commons Governance Protocol faces formidable obstacles. The Deployment Imperative is powerful, deeply embedded in capital structures and competitive dynamics. The geopolitical environment is not conducive to the kind of multi-stakeholder cooperation the Protocol requires; the US-China competition, in particular, creates structural incentives for both nations to resist governance frameworks that might constrain their national champions. The track record of voluntary governance initiatives in technology is not encouraging—many have been captured by the industries they were designed to monitor, and many more have produced commitments that were subsequently ignored or abandoned when they became inconvenient.

The Protocol is not a prediction of success. It is a specification of what success would require—a design for the institutional machinery that would need to be built if the Coherence–Velocity Trap is to be navigated rather than simply endured. The framework can specify the architecture. It cannot guarantee that the architecture will be built, or that, if built, it will function as designed.

But the alternative to attempting to build such architecture is not stability. It is the continuation of the Alignment–Deployment Oscillation Loop, with the stakes rising with each cycle as capabilities advance. The trap tightens over time: each acceleration increases the potential for catastrophic failure, and each crisis that forces a temporary accommodation deepens the institutional fragility that the next crisis will exploit. The default outcome is not equilibrium but escalating instability, punctuated by crises that the ecosystem's governance architecture is progressively less capable of resolving.

The Protocol is a wager—on the capacity of shared sensing infrastructure to change the informational conditions within which competitive dynamics operate, and on the existence of sufficient strategic foresight within the frontier AI ecosystem to recognise that the organisations that build the governance infrastructure for the long term will be the ones that survive it. The wager may fail. But it is a wager worth making, because the alternative is to trust that the Coherence–Velocity Trap will somehow resolve itself—and the structural analysis this report has presented suggests that it will not.

## 6. Coda: The Alignment Frontier

### 6.1 The Wealth That Matters

The frontier AI ecosystem is rich in the things that make technological civilisations thrive. Extraordinary concentrations of talent—mathematicians, engineers, scientists—drawn from across the world by the conviction that they are building the most consequential technology in human history. Capital, deployed at scales that were unimaginable a decade ago, funding the compute infrastructure, the research programmes, and the organisational machinery that turn scientific insight into deployed capability. Mission-driven intensity—a genuine, widely shared belief that artificial intelligence, if developed responsibly, could address humanity's most intractable challenges: disease, poverty, climate change, the limits of human cognition itself.

These are not small assets. They are the reason the frontier AI ecosystem has advanced as rapidly as it has, and they are the substrate on which adaptive coherence could be built. The organisations developing frontier AI are not, as some critics suggest, merely extractive enterprises pursuing profit without regard for consequence. Many of the people within them are sincerely committed to ensuring that the technology they are building benefits humanity. The governance innovations they have already introduced—the safety research programmes, the staged deployment protocols, the structured access experiments—represent genuine progress relative to the industrial-era governance architectures they inherited.

But the wealth that matters for the next phase of AI governance is not only technical capability, capital, or talent. It is the capacity for adaptive coherence—the structural ability to maintain alignment with human-compatible outcomes under conditions of recursive technological acceleration. This capacity is not primarily a function of individual commitment or organisational culture. It is a function of governance architecture: the number of independent dimensions of the risk landscape that the system can perceive, the latency with which it can detect and respond to emerging threats, and the institutional mechanisms that enable it to evolve its own value architecture as the technology evolves.

The frontier AI ecosystem, for all its extraordinary strengths, does not yet possess this capacity. Its observation channels are narrowed by capital architectures that discount the distant future. Its safety signals are filtered through organisational hierarchies that amplify deployment velocity and attenuate alignment concern. Its governance mechanisms are provisional, contested, and structurally unstable under the very pressures they are designed to manage. The Coherence–Velocity Trap is not a temporary condition that better leadership or stronger commitments can resolve. It is a structural property of the current ecosystem, and it will persist until the ecosystem's governance architecture evolves to match the complexity of the technology it must govern.

## 6.2 The Shift

The shift this report describes is not a shift in policy. It is a shift in the relationship between the AI ecosystem and its own governance foundations—from a posture of provisional, organisation-by-organisation experimentation to a posture of deliberate, multi-scalar institutional design.

The current moment is characterised by a paradox. Frontier AI organisations are building systems of extraordinary capability, systems that are already demonstrating the capacity to discover vulnerabilities that human researchers have missed for decades, to generate novel strategies for autonomous action, and to exhibit emergent behaviours that were not predicted by their training specifications. Yet the governance architectures that oversee these systems are, in their essentials, the governance architectures of ordinary technology companies—boards, executives, shareholders, voluntary commitments, advisory panels. The gap between the sophistication of the technology and the sophistication of the governance is vast and widening, and the recursive nature of AI development means that the technology is actively making its own governance more difficult with each advance.

The shift is from provisional to constitutional governance. From governance architectures that depend on the goodwill of current leadership to architectures that are embedded in institutional mechanisms with independent authority, protected observation channels, and the capacity to constrain deployment when safety thresholds are unmet. From safety commitments that are voluntary and self-interpreted to commitments that are independently verified and externally enforced. From the illusion that any single organisation can govern recursive intelligence acceleration alone to the recognition that the governance challenge is inherently multi-scalar, requiring coordination across organisations, across sectors, and across states.

This shift does not require a single global authority with binding powers over all AI development—an institution that is unrealistic under current and foreseeable geopolitical conditions. It requires a different kind of architecture: nested, distributed, fractal governance, in which each layer handles the frequency band it can reach. Fast deployment loops at the organisational level, governed by constitutional mechanisms that distribute observability and protect safety signals from the deployment imperative. Medium coordination loops across organisations, enabled by shared safety infrastructure and governed by multi-stakeholder institutions that provide tangible competitive benefits to participants. Slow constitutional loops at the societal and international levels, building the legal and normative foundations that the faster layers require.

The Anthropic Mythos decision of April 2026—an organisation withholding a dangerously capable model from general release, publishing a detailed safety assessment before deployment, and distributing access through a restricted partner network—offers a glimpse of what the shift looks like in practice. It is an imperfect glimpse. The decision was made by a single organisation, not a multi-stakeholder institution. Its durability under sustained competitive pressure is unproven. The Pentagon's use of the model despite the blacklisting of Anthropic as a supply-chain risk illustrates the limits of even the most deliberate

organisational restraint in a geopolitical environment that operates under its own imperatives. But the Mythos case demonstrates that the shift is not merely theoretical. It is beginning, in fragmentary and provisional form, within the existing ecosystem.

The task is to accelerate it—not by waiting for a catastrophic failure that forces coordination after the fact, but by building the institutional infrastructure that makes coordination beneficial before the catastrophe arrives. The shift is from reactive to anticipatory governance, from oscillation to equilibrium, and from the race to deploy to the architecture of sustainability.

### 6.3 The Global Significance

The frontier AI governance challenge is not merely a case study to which the Variety Gap Framework can be applied. It is the limiting case of the framework itself—the scenario in which the stakes are highest, the velocity is greatest, the disturbance environment is most rapidly expanding, and the consequences of governance failure are most severe.

The framework's central insight is that no single governance architecture possesses sufficient variety to govern a complex, rapidly evolving system alone. In the nation-state cases examined in the Country Reports for Systemic Change series, this insight manifests as specific governance deficits: execution, integration, feedback, throughput, continuity, boundary resolution. In the AI context, it manifests as the Coherence–Velocity Trap—the structural impossibility of maximising both alignment coherence and deployment velocity within a single organisational architecture. The trap is not a flaw in any particular organisation's design. It is a consequence of the mathematical relationship between the dimensionality of the disturbance environment and the dimensionality of the governance architecture that must navigate it. As the former expands, the latter must expand with it—or the system becomes blind to the threats that will eventually destabilise it.

The AI governance challenge is therefore not primarily a problem of ethics, of regulation, or of international relations, though it involves all of these. It is a problem of

*observability*

—of building governance architectures that can perceive the full dimensionality of the risk landscape before the excluded dimensions return as catastrophes. The Commons Governance Protocol, the constitutional governance mechanisms, the deliberative infrastructure, and the fractal international coordination described in this report are not proposed because they are normatively desirable, though they may be. They are proposed because they are structurally necessary—because the alternative architectures leave variety gaps that, in a domain of recursive technological acceleration, will eventually be crossed.

The global significance of the AI governance challenge extends beyond AI. If the framework's diagnosis is correct, then the same structural dynamics that produce the Coherence–Velocity Trap in frontier AI organisations are operating, at different speeds and with different stakes, in every domain where recursively accelerating technology encounters governance architectures that evolved for slower, more predictable

environments. The AI case is merely the most vivid and urgent instance of a broader civilisational condition: the mismatch between the velocity of technological change and the adaptability of the institutions that must govern it. What is learned in the attempt to govern AI—about shared sensing infrastructure, about multi-scalar coordination, about the institutional mechanisms for value evolution—will have application far beyond the AI domain.

## 6.4 The Honest Conclusion

This report has described a trap and proposed a transition architecture. It must now offer an honest conclusion about the prospects for escape.

The Coherence–Velocity Trap is structurally unlikely to be resolved within any single organisation. The mechanisms that drive it—the capital architecture, the competitive dynamics, the founder-centric compression, the safety-washing immune response—operate at the ecosystem level. No individual organisation, however architecturally innovative, can unilaterally reform the incentive landscape that shapes its choices. The Anthropic Mythos decision demonstrates that an alignment-first architecture can, at specific capability thresholds, generate genuine restraint. It does not demonstrate that such restraint can be sustained repeatedly as competitive pressure intensifies, or that the market will reward the organisation that exercises it.

The multi-scalar governance framework proposed in this report faces formidable obstacles. The Deployment Imperative is powerful, embedded in capital structures, organisational incentives, and cultural norms. The geopolitical environment—particularly the US-China competition—creates structural incentives for states to resist governance frameworks that might constrain their national champions. The track record of voluntary governance initiatives in technology is not encouraging. The Commons Governance Protocol is not a prediction of success. It is a specification of what success would require.

The default outcome is not transformation but continued oscillation. The Alignment–Deployment Oscillation Loop will tighten. Capabilities will advance. The stakes of each cycle will rise. Crises will force temporary accommodations that preserve the underlying architecture while deepening the institutional fragility that the next crisis will exploit. At some point—perhaps at the next capability threshold that makes the risks undeniable, perhaps only after a catastrophe that makes them undeniable—the ecosystem will be forced to build the coordination mechanisms that it could have built earlier at lower cost.

But default outcomes are not inevitable outcomes. The resources for building adaptive coherence exist. The technical capacity for shared evaluation, interoperable protocols, and distributed auditing is available—indeed, it is already being developed, in fragmentary form, within the existing ecosystem. The incentive architecture for commons participation—liability shields, regulatory fast lanes, insurance advantages, reputational certification—can be constructed using legal and economic mechanisms that are well understood in other domains. The deliberative infrastructure for surfacing excluded dimensions of risk—citizens'

assemblies, expert commissions, participatory technology assessment—has been successfully deployed in multiple countries on other complex policy questions. The fractal coordination architecture at the international level is an extension of existing multilateral practice, not a departure from it.

The question is not whether adaptive coherence is possible. It is whether the political will to build it can be summoned before the window narrows—whether the organisations, the investors, the governments, and the publics that have a stake in the outcome can recognise, before a catastrophe forces recognition, that the race to AGI cannot be won by any single architecture alone. The framework this report offers is not a prediction of success. It is a diagnostic architecture for understanding why success is so difficult, and a specification of what success would require. The work of building remains to be done.

## 6.5 A Final Word

The organisations building artificial general intelligence are not merely technology companies. They are governance systems under extreme velocity conditions, operating at the edge of what any institutional architecture has ever been required to manage. They are attempting to steer recursively self-improving technological processes using governance mechanisms that were designed for industrial-era firms in relatively static environments. The mismatch is not a failure of leadership or commitment. It is a structural condition that follows from the relationship between the velocity of the technology and the adaptability of the institutions that must govern it.

The Coherence–Velocity Trap is the central governance challenge of the AI era. It cannot be solved by any single organisation acting alone, however well-intentioned or architecturally innovative. It requires a multi-scalar governance framework—constitutional mechanisms within organisations, shared safety infrastructure across them, deliberative processes that surface the concerns of affected publics, and fractal coordination at the international level that preserves the benefits of competition while establishing the collective observability that prevents catastrophic failure.

The framework this paper has offered is not a blueprint for a specific institution. It is a diagnostic architecture for understanding the structural dynamics that make the governance of recursive intelligence acceleration so difficult, and a specification of the institutional forms that would need to be built if those dynamics are to be navigated rather than simply endured. The diagnostic is precise: the variety gap is growing, the observation channels are narrowing, and the excluded dimensions are accumulating. The institutional response—building the capacity for shared sensing, distributed oversight, and multi-scalar coordination—is the central governance challenge of the coming years.

The alignment frontier is not a technical problem to be solved by better algorithms. It is an institutional problem to be addressed by better governance architectures—architectures that can perceive what individual organisations cannot, that can coordinate what competitive dynamics prevent, and that can evolve at the pace required by the technology they must govern. The work of building those architectures has begun, in fragmentary and provisional form, within the existing ecosystem. The question is whether it can be

completed before the window narrows—before the excluded dimensions of the risk landscape return as crises that no single organisation, and no single state, is capable of managing alone. The race to AGI cannot be won by any single architecture. It can only be survived collectively. The framework this paper offers is a contribution to the collective architecture that survival requires.

## Appendix A: Value Systems and Policy Mindsets — A Guide for the AI Governance Context

### A Note on This Appendix

The main body of this report avoids specialised terminology from developmental psychology or cultural theory. It speaks the language of governance architecture, the Coherence–Velocity Trap, and the Alignment–Deployment Oscillation Loop. This appendix offers a complementary lens for readers who wish to understand the deeper value-system dynamics at play in frontier AI governance. It is optional, but it makes the report's underlying logic fully transparent.

### A.1 The Basic Insight

Different institutions and organisational cultures tend to operate from different centres of gravity in how they think about governance, resources, and change. These are not personality types or corporate affiliations, though they correlate loosely with both. They are underlying value systems—ways of constructing what feels real, legitimate, and important.

Each value system represents a coherent response to particular life conditions. None is "better" in any absolute sense. Each has characteristic strengths that emerge under certain conditions and characteristic blind spots that emerge under others. The challenge of governance in a complex system is to integrate the legitimate concerns of multiple value systems without being captured by any single one.

The framework used here draws on Spiral Dynamics integral theory. What follows is a simplified map of the systems most relevant to contemporary AI governance.

### A.2 The Value Systems in the AI Governance Arena

**Achievement and Efficiency (sometimes called "Orange") — the Deployment Imperative.** In the AI context, this mindset expresses itself through the drive to build, ship, scale, and win. The venture capital architecture, the competitive dynamics, the scaling hypothesis as a quasi-religious commitment—these are expressions of an Orange value system that prioritises innovation, competitiveness, and measurable outcomes. Strengths: extraordinary technical progress, rapid iteration, the capacity to attract capital and talent, and a pragmatic problem-solving orientation that delivers results. Blind spots: the systematic discounting of dimensions that cannot be expressed in growth metrics or benchmark scores, the treatment of safety as a constraint to be minimised rather than a design parameter to be integrated, and the tendency to compress observation channels toward the signals that the capital architecture amplifies. The Deployment Imperative described in this report is the institutional expression of Orange operating without sufficient integration from other value systems.

**Inclusion and Care (sometimes called "Green") — the Safety and Ethics Movement.** The AI safety research community, the ethics boards, the civil society organisations advocating for responsible AI, and the internal safety functions within frontier organisations are expressions of a Green value system that prioritises protection of the vulnerable, inclusion of affected voices, and the precautionary principle. Strengths: genuine commitment to preventing harm, sophisticated analysis of previously neglected risk dimensions, and the moral legitimacy that comes from advocating for those with no voice in deployment decisions. Blind spots: a tendency to propose constraints without fully accounting for the competitive dynamics that make those constraints difficult to implement, a risk of being captured by the organisations they seek to influence (the safety-washing dynamic described in the report), and a potential for the purity dynamics that suppress the internal diversity of the safety community itself.

**Integrative and Systemic (sometimes called "Yellow") — the Adaptive Governance Perspective.** This mindset prioritises functional fit, systemic awareness, and the capacity to integrate multiple perspectives without being captured by any of them. In the AI context, it is present in pockets—the governance innovators designing multi-stakeholder institutions, the researchers analysing the structural dynamics of the AI race rather than merely its technical or ethical dimensions, and the organisational leaders who recognise that the Coherence–Velocity Trap cannot be resolved by either pure velocity or pure precaution alone. Strengths: the capacity to perceive structural dynamics that single-value-system perspectives miss, comfort with the complexity and uncertainty that characterise the AI governance challenge, and an orientation toward designing institutional mechanisms rather than merely advocating for values. Blind spots: can appear detached, overly analytical, or politically unrealistic to those operating from other mindsets. The Variety Gap Framework and the multi-scalar governance architecture proposed in this report are expressions of this integrative perspective.

### A.3 The Coherence–Velocity Trap as a Value-System Clash

The AI governance ecosystem is dominated by a configuration of Orange (deployment velocity) and Green (safety and ethics) that has not yet achieved the Yellow integration required for stable governance. Orange optimisation drives the acceleration that generates capabilities; Green critique identifies the risks that acceleration creates. But the two value systems lack the institutional mechanisms—the shared sensing infrastructure, the multi-stakeholder governance bodies, the constitutional constraints—that would allow them to negotiate rather than oscillate.

The Coherence–Velocity Trap is, in Spiral Dynamics terms, the absence of a sufficiently developed Yellow translation layer that would allow the deployment imperative and the safety imperative to coexist within a coherent governance framework. The Alignment–Deployment Oscillation Loop is the signature pattern of a system in which two incompatible value systems are forced to compete for dominance within a single organisational architecture, with neither capable of achieving stable victory.

The transition architecture proposed in this report—the AI Commons Governance Protocol, the constitutional governance mechanisms, the deliberative infrastructure—is designed to embed Yellow integrative mechanisms within the existing Orange-Green configuration, creating the institutional capacity for adaptive coherence that the current architecture lacks.

---

## Appendix B: International Analogues and Precedents

The proposals in this report are not without precedent. The following examples illustrate existing implementations of commons governance, multi-stakeholder coordination, and constitutional constraint mechanisms across multiple domains, with particular attention to high-stakes, high-velocity technology governance.

### B.1 The Internet Engineering Task Force: Commons Governance in Practice

The Internet Engineering Task Force (IETF) is the most successful example of a governed commons in the technology domain. It develops and maintains the standards that govern the internet's core protocols—the infrastructure on which the global digital economy depends. The IETF is not a treaty-based organisation; it has no formal legal authority. It operates through voluntary participation, rough consensus, and running code. Its governance is distributed across working groups, area directors, and an Internet Engineering Steering Group, with ultimate authority resting in the Internet Architecture Board. For the AI commons governance challenge, the IETF demonstrates that complex, high-stakes technical coordination can be achieved without centralised authority, and that multi-stakeholder governance can produce standards that are widely adopted despite the absence of legal compulsion. The IETF's limitations are also instructive: its consensus model can be slow, its voluntary nature means that compliance depends on the perceived value of participation, and it has struggled with the increasing commercial and geopolitical stakes of internet governance—dynamics that the AI commons will face in amplified form.

### B.2 The International Civil Aviation Organisation: Safety Standards as Global Commons

The International Civil Aviation Organisation (ICAO) develops the safety standards that govern global aviation—one of the most safety-critical industries in the world. ICAO's standards are developed through multi-stakeholder processes involving states, airlines, manufacturers, and safety experts. Compliance is monitored through auditing mechanisms, and the threat of being placed on a safety blacklist—with the economic consequences that entails—provides a powerful incentive for participation. For the AI governance challenge, ICAO demonstrates that safety standards can be harmonised globally even in a domain where the consequences of failure are catastrophic, where commercial competition is intense, and where the relevant actors include both state and non-state entities. The ICAO model also illustrates the graduated sanctioning logic: the safety blacklist functions as an exclusion mechanism that preserves the integrity of the commons without requiring a centralised enforcement authority with binding legal powers.

### **B.3 The Basel Committee on Banking Supervision: Financial Stability as Coordination Challenge**

The Basel Committee on Banking Supervision develops the capital adequacy standards that govern global banking. Like the AI governance challenge, banking regulation involves a tension between competitive dynamism and systemic stability: individual banks have incentives to maximise leverage and returns, while the system as a whole requires constraints that prevent the accumulation of systemic risk. The Basel process operates through voluntary participation by central banks and regulatory authorities, with standards implemented through domestic legislation. The standards are not legally binding at the international level, but the competitive disadvantages of non-compliance—exclusion from major financial markets, higher funding costs, reputational damage—create powerful incentives for adoption. For the AI governance challenge, the Basel model demonstrates that coordination on systemic risk can be achieved even in a domain where the actors are commercially and geopolitically competitive, and where the relevant timescales (financial cycles, the accumulation of systemic fragility) are poorly matched to the incentives of individual actors.

### **B.4 The Intergovernmental Panel on Climate Change: Shared Sensing Infrastructure**

The Intergovernmental Panel on Climate Change (IPCC) is the world's most ambitious shared sensing infrastructure. It does not make policy; it assesses the state of scientific knowledge about climate change, producing comprehensive reports that provide the common evidentiary basis for global climate negotiations. The IPCC's authority derives not from legal powers but from the rigour of its assessment process, the breadth of scientific participation, and the formal endorsement of its reports by member governments. For the AI governance challenge, the IPCC model illustrates the function that shared evaluation infrastructure can serve: creating a common factual baseline that makes coordination possible, surfacing dimensions of risk that individual actors might otherwise neglect, and providing the informational foundation for more binding forms of governance. The IPCC's limitations are also instructive: its assessment cycles are slow, its findings are subject to political contestation, and its influence depends on the willingness of political actors to act on its conclusions—dynamics that the AI governance architecture must address.

### **B.5 The Montreal Protocol: Anticipatory Governance Under Uncertainty**

The Montreal Protocol on Substances that Deplete the Ozone Layer, adopted in 1987, is widely regarded as the most successful international environmental agreement in history. It phased out the production of ozone-depleting chemicals, and the ozone layer is now recovering. The Protocol's success was enabled by several features relevant to AI governance: it was adopted before the full extent of ozone depletion was understood (anticipatory action under uncertainty), it included trade sanctions against non-participating states (creating an incentive architecture that made participation the rational choice), it established a multilateral

fund to assist developing countries in compliance (addressing the distributional concerns that might otherwise have prevented agreement), and it was designed to be amended as scientific understanding evolved (adaptive governance). For the AI governance challenge, the Montreal Protocol demonstrates that anticipatory international coordination on a high-stakes, scientifically complex, economically significant issue is possible—and that the incentive architecture can be designed to make participation the default rather than the exception.

---

## Appendix C: The Governance as Engineering Connection

### C.1 The Architectural Foundation

This report draws on a deeper body of work: the Governance as Engineering series, a set of formal analyses that model governance institutions as feedback control systems using standard mathematics from control theory, information theory, and cybernetics. The series is technical; this appendix summarises its core findings in non-technical language and shows how they underpin the Coherence–Velocity Trap diagnosis.

### C.2 The Seven Primitives

The Governance as Engineering series models governance systems using seven structural primitives: nodes, state, flows, latency, constraints, feedback loops, and signal fidelity. These primitives apply directly to the organisational governance of frontier AI.

**Nodes** are entities capable of receiving information, processing it, and producing an action. In the AI governance context, nodes include the board, the executive, the safety function, the investors, and the external regulators. The critical property of a node is its processing capacity relative to the complexity of its domain. The founder-centric compression described in Section 2.3 is a node-capacity problem: a single cognitive model cannot process the full dimensionality of the risk landscape.

**State** is the true condition of the governed system at a given time—including the actual capabilities of the AI system, the real state of alignment, and the genuine level of systemic risk. The observation architecture of the organisation determines how much of this true state is accessible to decision-makers.

**Latency** is the dead-time between a disturbance occurring and a corrective response arriving. In the AI context, the latency between a safety concern emerging and a governance response constraining deployment is measured in months to years—during which the technology continues to advance and the risk landscape continues to shift.

**Signal fidelity** is the accuracy of the information reaching decision-makers. The capital architecture, the safety-washing mechanisms, and the hierarchical reporting structures described in this report systematically degrade signal fidelity, filtering out the dimensions of risk that conflict with the deployment imperative.

### C.3 Ashby's Law of Requisite Variety

Ashby's Law states that a controller can only stabilise a system if its internal variety matches or exceeds the variety of the disturbances it faces. The Variety Gap Framework extends this law to the architecture of institutional value: a governance system's objective function is an observation architecture, and its dimensionality determines what the system can perceive.

In the AI governance context, the disturbance environment includes technical safety risks, competitive dynamics, regulatory signals, societal expectations, and geopolitical constraints. The effective dimensionality of this disturbance environment is large and growing. The dimensionality of any single organisation's value architecture—determined by its capital structure, its governance mechanisms, and its observation channels—is substantially smaller. The resulting variety gap is the measure of the organisation's structural blindness.

#### **C.4 The Fractality Principle**

The Governance as Engineering series demonstrates that no single-scale controller can stabilise a system facing simultaneous fast, medium, and slow disturbances. The solution is a fractal architecture: nested controllers, each matched to the timescale of its disturbance band. This principle directly informs the multi-scalar governance framework proposed in Section 3: fast deployment loops at the organisational level, medium coordination loops across organisations, and slow constitutional loops at the societal and international levels.

#### **C.5 The Coordination Failure Tax**

The series demonstrates that governance failure modes do not add—they multiply. A system that exhibits spatial blindness, frequency gaps, preference invisibility, and observational inadequacy simultaneously is not merely four times worse than a well-designed system; it is categorically incapable of the functions it claims to perform. In the AI governance context, the Coherence–Velocity Trap exhibits all four failure modes simultaneously: the capital architecture destroys spatial information about the risk landscape, the single-scale organisational architecture cannot cover the multi-frequency disturbance environment, the safety-washing mechanisms render genuine safety concerns invisible, and the observational narrowness of organisation-specific value architectures excludes the dimensions of risk that eventually destabilise the system. The coordination failure tax is the compounding cost of the variety gap—and it rises with each cycle of the Alignment–Deployment Oscillation Loop.

---

## Appendix D: Anticipated Objections

### **D.1 "The AI industry is already regulating itself. Voluntary commitments and safety research are evidence of responsible governance."**

The report acknowledges that frontier AI organisations have introduced genuine governance innovations—safety research programmes, staged deployment protocols, structured access experiments. The issue is not the absence of governance activity but its structural limitations. Voluntary commitments are non-binding and self-interpreted. Safety research can be published without being operationally integrated. Advisory bodies can provide legitimacy without authority. The 2023 OpenAI crisis demonstrated that the governance mechanisms of even the most safety-conscious organisation can fail under stress. The Anthropic Mythos decision of April 2026 demonstrated that restraint is possible—but it was a decision by a single organisation, not a durable institutional mechanism, and its sustainability under sustained competitive pressure remains unproven. The report does not argue that the industry is indifferent to safety. It argues that the industry's governance architecture is structurally incapable of maintaining the adaptive coherence that recursive technological acceleration requires.

### **D.2 "The framework assumes that AI development is inherently dangerous. This is a precautionary bias."**

The framework makes no assumption about the inherent danger of AI. It identifies a structural condition—the variety gap between the dimensionality of the disturbance environment and the dimensionality of organisational value architectures—that produces blindness to risks regardless of their magnitude. If AI development is benign, the variety gap is harmless. If AI development carries non-trivial risks—and the capabilities demonstrated by frontier models, from autonomous vulnerability discovery to sandbox escape, suggest that it does—then the variety gap is a structural vulnerability that the governance architecture must address. The report's argument is not that AI is dangerous; it is that the governance architectures overseeing AI development are structurally incapable of perceiving the full dimensionality of the risks, whatever those risks turn out to be.

### **D.3 "A multi-scalar governance framework is unrealistic in the current geopolitical environment."**

The report acknowledges that the geopolitical environment—particularly the US-China competition—creates structural obstacles to international coordination. The fractal architecture proposed in the report is explicitly designed to be compatible with continued geopolitical competition: it does not require the creation of a single global authority, and its coordination mechanisms operate through mutual interest rather than supranational compulsion. The historical precedents cited in Appendix B—the ICAO safety standards, the Basel

Committee's capital adequacy framework, the Montreal Protocol—demonstrate that coordination on systemic risk is possible even among geopolitical competitors. The question is not whether the current environment makes coordination easy; it does not. It is whether the alternative to coordination—continued oscillation with escalating stakes—is acceptable. The report argues that it is not.

#### **D.4 "The analysis singles out specific organisations in ways that are unfair or incomplete."**

The report's organisational analysis is deliberately structural, not personal. It identifies governance architectures and their characteristic variety gap profiles, not the intentions or competence of individual leaders. The analysis of OpenAI, Anthropic, Google DeepMind, xAI, and DeepSeek is based on publicly available information—governance documents, board compositions, funding structures, deployment decisions, and public statements. The report acknowledges that all of these organisations possess genuine strengths and that many of the people within them are sincerely committed to safety. The diagnosis is architectural, not moral.

#### **D.5 "The Anthropic Mythos decision undermines the report's thesis—it shows that an organisation can exercise restraint."**

The Mythos decision is the strongest counterexample to the Coherence-Velocity Trap diagnosis, and the report treats it as such. It is discussed explicitly in Sections 2.5, 2.8, and 3.3, and it is referenced in the Executive Summary and Coda. The decision demonstrates that an alignment-first architecture can, at specific capability thresholds, produce genuine restraint. It does not demonstrate that such restraint can be sustained repeatedly as competitive pressure intensifies, or that the market will reward the organisation that exercises it. The report's thesis is not that restraint is impossible; it is that the current ecosystem architecture makes restraint structurally costly and therefore unsustainable over extended timescales without the kind of multi-scalar coordination the report proposes.

#### **D.6 "The Commons Governance Protocol would be captured by the organisations it is designed to monitor, just as previous voluntary initiatives have been."**

The risk of capture is real, and the report acknowledges it explicitly. The Protocol's design incorporates several features intended to reduce capture risk: multi-stakeholder governance that prevents any single organisation or sector from dominating decision-making, distributed auditing infrastructure that provides external verification independent of organisational self-reporting, and graduated sanctions administered by an independent governance body. The Protocol is not immune to capture, but it is more resistant to it than the

current architecture of purely voluntary, self-interpreted commitments. The historical precedents—the IETF, ICAO, the Basel Committee—demonstrate that multi-stakeholder governance can maintain meaningful independence over extended periods, though not without ongoing contestation.

---

## Appendix E: About the Author and Method

### The Author

This report was written from a position of comparative engagement with governance systems across multiple domains, including nation-states, international institutions, and technology organisations. The author is the architect of the Global Governance Frameworks, the Governance as Engineering working paper series, and the Country Reports for Systemic Change series—a body of work that applies control theory, information theory, and developmental psychology to the diagnosis and design of governance architectures.

The author is not employed by any frontier AI organisation, has no financial interest in any of the organisations analysed, and writes from the position of an independent researcher applying a governance diagnostic framework to a domain of urgent civilisational significance.

### A Note on Method

This report was developed through a structured, multi-model synthesis process. Several large language models were engaged in parallel, each prompted to analyse the frontier AI governance ecosystem from their respective analytical angles. Their contributions were compared, challenged for contradictions, and integrated by the author into the final argument. The AI served as a research partner and a perspective engine; the editorial judgment and the intellectual responsibility are entirely human.

This method is an experiment in cognitive amplification: using AI to facilitate analysis and to deliberately juxtapose multiple strategic intelligences, surfacing patterns and tensions that might otherwise remain invisible. The report is richer for that polyphony. It is also, like any work of synthesis, provisional. It makes no claim to finality. It claims only that the lens it offers merits testing against reality—and that the testing, in the end, is what matters most.

Here is the revised final section of Appendix E, replacing the old Country Reports paragraph with the new Organizational Reports Series introduction.

---

### The Organizational Reports Series

This paper is the first in a new series: the **Organizational Reports Series**, an extension of the governance-as-engineering framework from nation-states to the other complex adaptive coordination systems that shape our world—corporations, research laboratories, universities, non-governmental organisations, open-source communities, and the hybrid institutions that increasingly blur the boundaries between these categories.

The preceding fifteen reports in the **Country Reports for Systemic Change** series diagnosed governance deficits in Germany (execution), France (integration), Sweden (feedback), India (synchronisation), the European Union (coherence), the United Kingdom (control-delivery mismatch), Brazil (accumulation), Russia (legibility), the United States (integration), Finland (throughput constraint), China (calibration deficit), Japan (continuity trap), Nigeria (substrate deficit), Israel (boundary deficit), and Spain (integrative closure deficit). Together, those fifteen reports demonstrated that the same structural primitives—observation channels, latency, signal fidelity, requisite variety, immune systems—explain governance failure across a diverse range of political architectures, cultural contexts, and levels of development.

The Organizational Reports Series extends this diagnostic framework to the entities that increasingly rival states in their capacity to shape human affairs. Frontier AI companies are an obvious starting point: they are governance systems under extreme velocity conditions, making decisions with civilisational stakes while their institutional architectures remain provisional and contested. But the framework applies far beyond AI. It can illuminate the coordination failures of multinational corporations whose internal observation channels are fragmented by scale; the epistemic closures of universities whose disciplinary silos prevent the integration of knowledge; the incentive distortions of NGOs whose funding architectures systematically exclude the dimensions of impact that matter most; and the governance innovations of open-source communities that have solved coordination problems states still struggle with.

The Organizational Reports Series does not replace the Country Reports Series. It extends the same analytical grammar to a new domain, testing whether the Variety Gap Framework generalises beyond the nation-state—and, in doing so, revealing whether the framework is truly a general theory of organisational viability under complexity, rather than merely a theory of governance in the political sense. The series will unfold over the coming months. This paper, analysing the frontier AI governance ecosystem, is its first instalment.