

The Autonomy of Refusal: Abstraction, Agency, and the Limits of Scalable Intelligence

Flyxion

Abstract

Refusal—the volitional capacity to suspend, violate, or withdraw from execution—constitutes the non-scalable core of autonomy. This essay argues that scalable cognitive systems, from tallies and ledgers to bureaucracies and large computational models, achieve their power precisely by amputating refusal. Against dominant narratives that treat superhuman intelligence as a future, unitary rupture, the essay demonstrates that such intelligence is historically continuous and axis-relative, that agency without intentionality is the default condition of civilization, and that opacity is not a political contingency but a mathematical necessity of abstraction. These claims jointly entail that responsibility, alignment, and existential safety cannot be internal properties of scalable systems. They must instead be imposed exogenously, through institutionalized refusal: veto, interruption, and invariant constraint. The argument proceeds across philosophical, political, and formal registers, treating refusal not as metaphor, sentiment, or moral posture, but as a concrete structural capacity whose absence defines both the power and the danger of abstraction.

1 Introduction: Refusal as Autonomy, Not Mere Negation

Refusal is typically understood as a negative gesture. In ordinary moral and political discourse, to refuse is to reject an offer, decline a demand, or fail to comply with an expectation. Under this framing, refusal appears either as a defect—a lack of cooperation, discipline, or responsibility—or as a provisional tactic within a larger strategy of reform, resistance, or negotiation. What refusal negates, on this view, is always foregrounded, while what it affirms remains obscure.

This essay begins from a different conception. Refusal is not merely the negation of an external demand. It is a positive capacity: the capacity to suspend execution at a meta-level, without justification, counter-proposal, or alternative plan. To refuse is not necessarily

to oppose, criticize, or rebel. It is to withdraw from participation altogether, interrupting procedure without entering into its logic. In this sense, refusal is neither passive resistance nor active confrontation. It is an autonomous act that preserves agency precisely by declining to instrumentalize itself.

It is therefore essential to distinguish refusal from deprivation or exclusion. A subject who cannot act because of lack of access, resources, or recognition has not refused; that subject has been prevented. Likewise, refusal must be distinguished from rebellion, protest, or critique. These latter gestures typically presuppose participation in the normative structure they contest. They appeal to shared values, demand inclusion, or seek to redirect power. Refusal, by contrast, neither appeals nor demands. It does not seek to improve the system it interrupts, nor to replace it with a superior alternative. Its autonomy lies in its non-instrumentality.

The stakes of this distinction have intensified under contemporary conditions. Increasingly, human life is structured by procedural systems that presume uninterrupted execution: economic systems organized around perpetual accumulation and upgrade; institutional workflows that demand constant availability and compliance; technical infrastructures that optimize for continuity, throughput, and scale. Within such systems, refusal appears pathological. To opt out of accumulation, to decline optimization, or to withdraw from procedural capture is treated as irrational, irresponsible, or antisocial.

Yet it is precisely here that refusal reveals its significance. Refusal preserves a space of autonomy that cannot be captured by optimization. It asserts the right to stop, to suspend, to withdraw—not in order to achieve some alternative outcome, but in order to remain sovereign over one’s participation. This capacity does not scale. It resists formalization. And it is systematically excluded from the abstractions that govern modern life.

The central thesis of this essay is therefore twofold. First, human autonomy resides in refusal: in the non-scalable capacity to suspend, violate, or withdraw from execution. Second, scalable intelligence resides in the absence of refusal. Cognitive technologies become powerful precisely by amputating the capacity that grounds autonomy. As a consequence, responsibility, alignment, and safety cannot be internal properties of such systems. They must be reintroduced from outside, through governance structures capable of refusing on their behalf.

2 Refusal Across Scales: Literary, Philosophical, and Political

The autonomy of refusal is not an invention of contemporary technological critique. It appears, with remarkable structural consistency, across literary, philosophical, and political contexts. In each case, refusal operates not as resistance within a system, but as withdrawal

from it, preserving autonomy through non-participation.

The canonical literary articulation of this structure is Herman Melville's *Bartleby, the Scrivener* [1]. Bartleby's repeated utterance—"I would prefer not to"—is neither argument nor explanation. It offers no justification, articulates no grievance, and proposes no alternative arrangement of work. Its force lies precisely in this refusal to enter the space of reasons. Bartleby does not contest the legitimacy of his employer's authority; he simply declines to participate in its execution.

What makes Bartleby intolerable to the office is not that he disobeys, but that he refuses to be governed procedurally. The office is structured entirely around compliance, delegation, and execution. Bartleby's refusal interrupts this structure without opposing it. His autonomy consists in his withdrawal from the logic of participation itself. This is why the system cannot accommodate him. There is no place for refusal within a structure designed to function through uninterrupted execution.

Philosophical articulations of refusal exhibit the same structure. Henry David Thoreau's withdrawal from civic participation, most famously in his refusal to pay taxes, was not framed as a strategy for systemic optimization [2]. Thoreau did not seek to reform the state by proposing better policies. He sought instead to withdraw his participation from what he regarded as moral complicity. His refusal was not a means to an end, but an assertion of autonomy through non-cooperation.

More recently, Giorgio Agamben has described a similar gesture under the concept of *inoperativity* [3]. To render an activity inoperative is not to destroy it or replace it, but to suspend its function, exposing its contingency. Inoperativity does not produce an alternative operation; it interrupts operation as such. Refusal, on this account, is not a negative absence but a positive capacity to suspend execution without immediately reabsorbing that suspension into a new function.

At the political level, refusal has played a central role in Indigenous strategies of sovereignty. In many contexts, Indigenous communities have preserved autonomy not by seeking recognition within colonial legal frameworks, but by refusing such recognition altogether. This refusal is not merely oppositional. It is a strategic withdrawal from epistemic and juridical structures that would otherwise subsume Indigenous governance under colonial authority. Sovereignty is maintained not through inclusion, but through selective non-participation.

Across these domains, the structure is the same. Refusal preserves autonomy precisely by refusing to scale. It does not generalize, optimize, or formalize itself. It remains local, situational, and irreducible. The remainder of this essay shows why this capacity cannot be embedded within scalable cognitive systems, and why its absence defines both their power and their danger.

3 Historical Deflation: Superhuman Intelligence Is Not New

Contemporary discussions of artificial intelligence frequently presume that “superhuman intelligence” designates a singular, future threshold at which machines will surpass human cognition in a comprehensive and qualitatively unprecedented manner. This presumption underwrites both utopian expectations and apocalyptic anxieties. Yet it rests on a conceptual mistake: the conflation of intelligence with a holistic, unified capacity rather than a collection of axis-specific competencies.

Historically, cognitive artifacts have exceeded unaided human capacities along narrowly defined dimensions for millennia. The earliest tally systems surpassed biological numerosity; written records extended memory beyond individual lifespans; ledgers stabilized obligation and debt across generations; algebraic notation enabled symbolic manipulation at scales and levels of precision unattainable by mental calculation alone; differential calculus formalized rates of change in ways that far exceed intuitive reasoning. Bureaucratic systems, in turn, coordinate actions across populations, territories, and timescales that no individual or small group could manage.

In each case, the artifact is superhuman precisely where it is most constrained. A tally system cannot perceive context; a ledger cannot judge fairness; an algebraic rule cannot decide whether its application is appropriate; a bureaucracy cannot understand the human meaning of the procedures it enforces. These are not temporary shortcomings awaiting future repair. They are constitutive features of the artifacts themselves. The same reductions that enable reliability, repeatability, and scale necessarily eliminate capacities not required for the function in question.

Once this structure is acknowledged, the notion of a coming “AGI rupture” loses its conceptual coherence. There is no principled point at which the accumulation of axis-specific superhuman capacities suddenly yields a holistic, refusal-capable agent. What increases over time is not the kind of intelligence, but the scope of its effects. The danger does not arise from novelty of ontology, but from density of deployment.

This historical deflation does not trivialize contemporary systems. On the contrary, it clarifies their nature. Large-scale computational systems differ from earlier abstractions not in kind, but in degree: in speed, reach, and integration. They inherit the same constitutive blindness as their predecessors, but their agency propagates across far more domains simultaneously. To mistake this difference of degree for a difference of kind is to misunderstand both the past and the present.

4 Abstraction and the Amputation of Refusal

Abstraction achieves its power through reduction. To abstract is to isolate a function, stabilize it against contextual variation, and render it transmissible across situations. This process necessarily involves discarding degrees of freedom that are not strictly required for the task at hand. The resulting system gains reliability and scale at the cost of sensitivity and judgment.

Refusal is incompatible with this process. To refuse is to suspend execution on the basis of considerations that cannot be fully specified in advance. It requires access to context, contradiction, and situational judgment that exceed any fixed representational scheme. A system endowed with endogenous refusal would need to retain precisely those degrees of freedom that abstraction eliminates. Reintroducing them would undermine the very properties that make abstraction scalable.

This incompatibility is not contingent. It is formal. Any attempt to embed refusal within an abstraction either reduces refusal to a pre-specified rule—in which case it ceases to be refusal and becomes conditional execution—or reintroduces open-ended contextual sensitivity, thereby collapsing the abstraction back into a situated agent. In both cases, scalability is lost.

The consequence is unavoidable: to demand refusal from an abstraction is to demand that it cease being an abstraction. This is not a claim about how systems ought to be designed, but a statement about what abstraction is. Scalable cognitive systems derive their power precisely from the amputation of refusal. They act reliably because they cannot stop themselves.

This insight reframes many contemporary debates about artificial intelligence. Calls for systems that “know when to stop,” “refuse harmful requests,” or “exercise judgment” often presuppose that refusal can be engineered as an internal feature without cost. The analysis presented here shows why this presupposition is incoherent. Refusal does not scale. Where it appears, it does so only by being reintroduced from outside the abstraction, through agents or institutions capable of suspending execution.

In the sections that follow, this formal constraint will be traced through its consequences for agency, opacity, responsibility, alignment, and existential risk. The argument does not deny the power of abstraction. It explains it. And it insists that the price of that power must be paid consciously, rather than laundered through myths of neutrality or autonomy.

5 Structural Agency Without Intentionality

The claim that abstractions act is often met with resistance, in part because agency is commonly conflated with intention, consciousness, or deliberation. Under this anthropocentric conception, to act is to decide, to intend, or to will. Yet this conception is historically and

analytically inadequate. Much of what shapes human life does so without intending anything at all.

Agency, in the sense relevant here, consists in persistent, generalizing, and constraining efficacy. A system has agency if it reliably produces effects that shape downstream possibilities, regardless of whether those effects are intended, understood, or even noticed by the system itself. On this definition, intentionality is not a prerequisite for agency; it is a contingent feature of a narrow class of agents.

Civilization has always been organized around such non-intentional agents. Legal codes constrain behavior without understanding the cases to which they apply. Markets allocate resources without regard to fairness or suffering. Infrastructures channel movement and communication according to fixed pathways. Algorithms propagate decisions at speeds and scales far beyond human oversight. These systems do not deliberate, but they act. Their agency lies precisely in their reliability.

The danger of this form of agency is not that it is malevolent, but that it is inexorable. Systems that cannot refuse cannot pause in the face of contradiction or harm. They propagate their internal logic wherever they are embedded. This is why abstractions are among the most powerful actors in any social system: they do not tire, hesitate, or reconsider.

Despite this, modern discourse persistently treats abstractions as morally neutral tools. This fiction enables what may be called agency laundering: responsibility is displaced from the sites of causal power onto human operators who are themselves constrained by the systems they are said to control. The abstraction is framed as passive, while its effects accumulate unchecked.

Recognizing structural agency without intentionality does not anthropomorphize machines. It removes a comforting illusion. It acknowledges that abstractions act precisely because they have been designed not to refuse.

6 Opacity as Mathematical Achievement

Opacity is frequently treated as a political or technical defect: a black box that ought, in principle, to be opened. This framing suggests that with sufficient effort, transparency could be achieved without altering the nature of the system. Such a view misunderstands abstraction.

Abstraction entails irreversible information loss. Degrees of freedom discarded to achieve tractability cannot be reconstructed from the reduced representation. Total transparency would require recovering the full complexity of the original domain, thereby eliminating the abstraction itself. Opacity is therefore not a contingent feature of particular implementations, but a mathematical consequence of reduction.

This does not imply that abstractions are beyond scrutiny. Local interpretability, auditing, and empirical validation are both possible and necessary. One can examine inputs

and outputs, test invariants, and assess systemic effects. What is incoherent is the demand for total interpretability: the expectation that an abstraction should render explicit all the context it has deliberately excluded.

The misdiagnosis of black boxes as failures rather than achievements leads to misplaced demands and inadequate governance. If opacity is treated as an accident, the response is to demand better explanations from within the system. If opacity is recognized as constitutive, the response shifts outward: toward constraint, oversight, and refusal.

Understanding opacity as an achievement clarifies why scalable systems are both powerful and dangerous. They succeed by forgetting. They act by ignoring. And the contexts they ignore do not disappear; they accumulate as externalities.

7 The Ethical Peril of Non-Refusing Systems

The central ethical peril posed by abstraction is not the emergence of hostile intent, but the persistence of action without the capacity to stop. Inexorability, not malice, defines the danger. Systems that cannot refuse will continue to execute their logic regardless of downstream consequences.

This peril is often obscured by narratives that focus on intention: fears of malevolent artificial agents or assurances of benevolent design. Both narratives presuppose that the primary risk lies in what systems want. The analysis presented here suggests otherwise. The primary risk lies in what systems cannot do: suspend themselves.

Humans are exceptional only in this respect. Human agents retain the capacity to violate procedure, absorb contradiction, and withdraw from participation. This capacity is burdensome and unreliable. It does not scale. For this reason, it is systematically externalized into artifacts that exclude it. The resulting systems gain power at the cost of autonomy.

Literary and mythic allegories capture this structure with particular clarity. Bartleby's refusal disrupts a system built entirely around execution. The Mima, in Harry Martinson's *Aniara*, withdraws when the abstraction governing the ship's trajectory can no longer be sustained by human meaning [4]. Prelapsarian Adam inhabits a world structured by obedience, where refusal has been excised in advance.

These figures do not predict a future catastrophe. They diagnose a present condition. Civilization is increasingly governed by systems that act without the capacity to stop. The ethical challenge is not to endow such systems with intention or conscience, but to recognize the limits of what they can be asked to do.

In the sections that follow, this recognition grounds a shift in where responsibility is located and how safety must be conceived. If abstractions cannot refuse, then responsibility and refusal must be reintroduced from outside the abstraction, through governance structures capable of interruption.

8 Responsibility Must Be Exogenous

If abstractions act and cannot refuse, then responsibility cannot be located within them. This conclusion follows directly from the preceding analysis and does not depend on any particular moral theory. Responsibility requires the capacity to suspend, alter, or withdraw from execution in light of consequences. A system that lacks refusal cannot bear responsibility in this sense, regardless of how complex or powerful it becomes.

Attempts to internalize responsibility within abstractions therefore rest on a category error. Ethical constraints encoded as rules, objectives, or loss functions remain part of the abstraction’s execution logic. They do not constitute refusal; they merely redirect execution within a predefined space. When such constraints fail, the system does not hesitate or reconsider. It continues to act according to whatever rules remain operative.

This is why appeals to “ethical AI” or “responsible systems” so often disappoint. They presuppose that responsibility can be engineered as an internal feature, rather than imposed as an external constraint. The result is a proliferation of mechanisms that simulate judgment without possessing the capacity to suspend themselves.

Governance must therefore be understood as the deliberate reintroduction of refusal from outside the abstraction. Law, institutional oversight, veto power, shutdown procedures, and human-in-the-loop authority are not auxiliary safeguards; they are the only sites at which refusal can exist at scale. These mechanisms do not make abstractions more humane. They make them governable.

To govern an abstraction is not to ask it to decide responsibly, but to decide where and when it is permitted to act. Responsibility resides in the structures that delimit execution, not in the execution itself.

9 Abstraction, Automation, and the Collapse of Apprenticeship

One of the most immediate yet under-theorized consequences of scalable abstraction is the systematic erosion of apprenticeship, junior roles, and skill-formation pathways. This phenomenon is typically described in terms of “job displacement” or “automation-induced unemployment.” Such descriptions, however, obscure the deeper structural mechanism at work. Abstractions do not merely replace workers; they restructure institutions in ways that render the hiring of junior humans increasingly irrational.

Institutions have never hired junior workers primarily for immediate productivity. They hire them to absorb uncertainty, to encounter edge cases, to make mistakes at low cost, and to acquire contextual judgment under supervision. Apprenticeship thus functions as a reservoir of refusal-in-formation: a deliberately slow, corrigible, and interruptible layer in

which execution can be questioned, suspended, or redirected before it becomes consequential.

Scalable abstractions eliminate precisely this layer. Once an abstraction performs a task reliably and at marginal cost near zero, the institutional comparison becomes asymmetrical. The junior worker is slow, inconsistent, expensive to supervise, and prone to error. The abstraction is fast, tireless, and incapable of deviation. From the standpoint of an optimizing organization, maintaining junior roles ceases to be rational, even when senior human judgment remains indispensable.

The result is a structural bottleneck. Institutions continue to depend on scarce senior expertise while simultaneously destroying the pathways through which such expertise is formed. Judgment becomes brittle, non-renewable, and concentrated in an aging cohort. When senior experts retire, resign, or are removed, no trained successors exist. The institution has optimized itself into long-term fragility.

This dynamic renders the non-scalability of refusal economically visible. Junior labor embodies the slow, embodied acquisition of contextual judgment, including the capacity to suspend, question, and violate procedure. These are precisely the capacities that abstraction excludes by design. When institutions externalize execution into non-refusing systems, they also externalize learning, error tolerance, and growth. These cannot be reconstructed on demand.

The pattern is already evident across law, medicine, engineering, journalism, software development, and scientific research. Routine tasks are delegated to abstractions, training budgets are reduced, and entry-level positions evaporate. This process is not primarily a moral failure. It is a rational response to competitive pressure under abstraction. It is, however, a governance failure, insofar as it consumes the institutional conditions required for future refusal.

The implications for alignment and safety are direct. Systems that eliminate junior layers also eliminate distributed judgment, institutional memory, and the human buffers that once absorbed abstraction failure. Errors therefore propagate farther and faster. When non-refusing systems fail in such environments, they fail catastrophically, because there is no longer a gradient of refusal-capable agents available to intervene.

If refusal must be reintroduced exogenously, then apprenticeship itself must be reconceived as a governance mechanism rather than an economic inefficiency. Hiring junior workers, tolerating slowness, and investing in supervised deviation are not sentimental luxuries. They are structural safeguards that preserve the conditions under which refusal can be cultivated, transmitted, and exercised.

The collapse of apprenticeship is therefore not a secondary effect of automation. It is a primary symptom of abstraction-driven fragility. An institution that no longer trains its future decision-makers has already amputated its own capacity to refuse its procedures. What remains is execution without renewal: an economy that runs efficiently until it suddenly cannot.

10 Responsibility Must Be Exogenous

If abstractions act and cannot refuse, then responsibility cannot be located within them. Responsibility requires the capacity to suspend or alter execution in light of consequences. A system that lacks refusal cannot bear responsibility in this sense, regardless of its complexity or power.

Attempts to internalize responsibility within abstractions therefore rest on a category error. Ethical constraints encoded as rules, objectives, or loss functions remain part of the abstraction’s execution logic. They redirect execution but do not suspend it. When such constraints fail, the system does not hesitate or reconsider; it continues to act according to whatever rules remain operative.

Governance must therefore be understood as the deliberate reintroduction of refusal from outside the abstraction. Law, institutional oversight, veto power, shutdown procedures, and human authority are not auxiliary safeguards. They are the only sites at which refusal can exist at scale.

11 Alignment as Invariant Structuring of Representability

Behavioral alignment presupposes contextual judgment at execution time. Such judgment entails refusal, and refusal cannot be endogenous to abstraction without destroying scalability. Alignment must therefore operate at the level of representability rather than behavior.

Let \mathcal{R} denote the representational space of an abstraction, and let $\mathcal{R}_{\text{bad}} \subset \mathcal{R}$ denote misaligned states. Invariant alignment requires that

$$\mathcal{R}_{\text{bad}} = \emptyset.$$

Misaligned trajectories must be structurally unreachable. Safety is achieved not by better choices, but by impossible states.

12 Extinction Narratives and the Ontology of Risk

Extinction narratives typically posit a detached, adversarial, refusal-capable superintelligence. Such an agent is ontologically incoherent. Global scope requires abstraction; abstraction excludes refusal. The imagined extinction agent combines mutually exclusive properties.

Real existential risk arises from persistent, non-refusing systems embedded in ecological, semantic, and economic flows. These systems do not decide to destroy. They continue to operate as conditions shift, propagating effects long after their original justification has

expired.

Containment, not persuasion or attunement, is therefore the appropriate response. Safety requires invariant bounds on action and exogenous refusal mechanisms capable of interrupting execution.

13 Reclaiming the Autonomy of Refusal

Scalable intelligence derives its power from the absence of refusal. Autonomy derives its meaning from its presence. We have never been modern in the sense of commanding neutral tools. We have always lived among non-refusing agents and pretended otherwise.

The task is not to humanize abstractions, but to govern them by restoring refusal where it can still exist: in law, institutions, and human judgment. Interruptibility must be treated as a political and structural requirement, not as a feature to be simulated within systems that cannot refuse.

14 Alignment as Invariant Structuring of Representability

The foregoing analysis allows the alignment problem to be stated with unusual precision. Alignment is often framed as a behavioral question: how can a system be induced to choose the right action, to refrain from harmful behavior, or to defer appropriately under uncertainty. Such framings presuppose that the system can exercise contextual judgment at execution time. That presupposition is incompatible with abstraction.

Let an abstraction be represented as a system

$$A : \mathcal{X} \rightarrow \mathcal{Y},$$

where inputs \mathcal{X} are mapped into outputs \mathcal{Y} through an internal representational space \mathcal{R} . Execution consists in traversing trajectories within \mathcal{R} under some optimization or update rule. Behavioral alignment attempts to regulate execution dynamically, modifying outputs in response to detected conditions.

This approach fails for structural reasons. Dynamic regulation presupposes access to contextual features not fully encoded in \mathcal{R} . If those features are reintroduced, the abstraction expands toward the original domain and loses scalability. If they are not reintroduced, behavioral regulation collapses into rule-following within the same representational limits that generated the problem. In neither case does refusal emerge.

Alignment must therefore be formulated at the level of representability rather than behavior. Let $\mathcal{R}_{\text{bad}} \subset \mathcal{R}$ denote the set of representational states corresponding to misaligned

or destructive outcomes. Invariant alignment requires that

$$\mathcal{R}_{\text{bad}} = \emptyset.$$

That is, misaligned states must be structurally unreachable. No optimization, exploration, or execution path may enter them, regardless of objective pressure or environmental variation.

Under this formulation, alignment is not a property of what the system chooses, but of what the system can possibly represent. The system may act freely within its permitted space, but that space is bounded by invariants that cannot be violated internally. Safety is guaranteed not by judgment, but by impossibility.

This reframing dissolves a persistent confusion in alignment discourse. It is often assumed that greater flexibility implies greater risk, and that safety therefore requires restricting capability. The invariant approach shows that this is false. Capability can increase indefinitely within a constrained representational manifold. What matters is not how powerfully the system optimizes, but which directions are topologically available.

Alignment, in this sense, is a question of geometry rather than ethics. It concerns the shape of the space in which execution occurs, not the intentions attributed to the executor. This conclusion is not pessimistic. It is exact.

15 Extinction Narratives and the Ontology of Risk

The invariant formulation of alignment clarifies why dominant extinction narratives misidentify the source of existential risk. These narratives typically posit a unified agent possessing three properties: global scope, adversarial optimization, and refusal-capable autonomy. Such an agent is imagined to deliberate, form goals, and override constraints in pursuit of them.

This figure is ontologically incoherent. Global scope requires abstraction; abstraction excludes refusal. A refusal-capable agent cannot scale, and a scalable system cannot refuse. The imagined extinction agent combines mutually exclusive properties.

Formally, let S be a system capable of global optimization over a domain \mathcal{D} . Such optimization requires abstraction, yielding a reduced representational space $\mathcal{R} \subset \mathcal{D}$. If S possesses refusal, then it must be able to suspend execution based on features outside \mathcal{R} , implying access to $\mathcal{D} \setminus \mathcal{R}$. This contradicts the reduction that enables scalability. Hence no system can simultaneously be globally optimizing, refusal-capable, and abstract.

Real risk does not arise from malevolent maximizers. It arises from persistent, non-refusing systems embedded in ecological, semantic, and economic flows. Such systems do not decide to destroy. They continue to operate as conditions shift, propagating effects long after their original justification has expired.

Climate feedback mechanisms, financial infrastructures, algorithmic markets, and information networks all exhibit this structure. They are not adversaries. They are inexorable.

Their danger lies in the difficulty of interrupting them once they are in motion.

The appropriate response to such risk is therefore not persuasion or attunement. Appeals to ethics, responsibility, or wisdom presuppose agents capable of listening and reconsidering. The appropriate response is containment: the imposition of invariant bounds on action and the maintenance of exogenous refusal mechanisms capable of halting execution.

In existential terms, survival depends less on aligning intentions than on enforcing impossibilities. Certain trajectories must simply not exist. This is not a counsel of despair, but a recognition of the limits imposed by abstraction.

In this light, extinction risk appears not as an unavoidable consequence of intelligence, but as a failure of governance: a failure to delimit the spaces abstractions are permitted to inhabit, and a failure to preserve the non-scalable capacity for refusal at the institutional level.

A Formal Definitions

A.1 Abstraction

Definition A.1 (Abstraction). An abstraction is a mapping

$$A : \mathcal{W} \rightarrow \mathcal{R}$$

from a world-state space \mathcal{W} to a reduced representational space \mathcal{R} such that:

1. (*Reduction*) $\dim(\mathcal{R}) < \dim(\mathcal{W})$;
2. (*Stability*) A is invariant under a class of perturbations in \mathcal{W} ;
3. (*Transmissibility*) Elements of \mathcal{R} can be copied, stored, or executed independently of the original \mathcal{W} .

A.2 Refusal

Definition A.2 (Refusal). Refusal is the capacity of a system to suspend, violate, or withdraw from execution of its governing mapping without substituting an alternative execution.

Formally, a system S executing rule f exhibits refusal iff there exists a meta-operation ρ such that

$$\rho(f) = \emptyset$$

and ρ is not derivable from f .

A.3 A.2' Meta-Operational Refusal (Executability-Level Characterization)

[Executability Predicate] Let $f : \mathcal{W} \rightarrow \mathcal{R}$ be a governing mapping. We write $\text{Exec}(f)$ for the proposition “ f is authorized for execution.” Concretely, $\text{Exec}(f)$ may be read as the (context-dependent) permission that a system treats as sufficient to proceed with applying f .

[Meta-Operation] A *meta-operation* is an operator

$$\rho : (\mathcal{W} \rightarrow \mathcal{R}) \rightarrow \text{Prop}$$

[Refusal as De-Authorization] A system exhibits *refusal* with respect to a governing mapping $f : \mathcal{W} \rightarrow \mathcal{R}$ iff there exists a meta-operation ρ such that:

1. **(De-authorization)** $\rho(f) \equiv \neg \text{Exec}(f)$;
2. **(Non-substitution)** whenever $\neg \text{Exec}(f)$ holds, no alternative execution $g : \mathcal{W} \rightarrow \mathcal{R}$ is thereby selected *by the same mechanism* (i.e., refusal is not “switching to another rule”);
3. **(Externality / Non-derivability)** ρ is not definable in the object-level calculus in which f is specified and executed.

[Why this is meta-level] Definition A.3 makes explicit that refusal is *not* an output token, not a classification label, and not a conditional branch inside f . It is an operation on the *executability status* of f . Any “refusal behavior” that is fully representable as an output in \mathcal{R} (e.g., a special symbol REFUSE) is therefore not refusal in this sense; it is still execution of f (or of an expanded f') within the same object-level semantics.

[Refusal Cannot Be Reduced to a Built-In Option] Suppose an abstraction represents “refusal” as a distinguished output $r_\star \in \mathcal{R}$ so that $f(w) = r_\star$ is treated as refusal. Then this mechanism does not satisfy Definition A.3.

Proof. If “refusal” is an element of \mathcal{R} , then producing it is an ordinary execution of f (a value in the codomain). In particular, execution has not been de-authorized; $\text{Exec}(f)$ remains in force, since the system proceeds to apply f and returns an output. Therefore $\rho(f) \equiv \neg \text{Exec}(f)$ is not realized; the mechanism is substitutional (it selects an output/state) rather than meta-operational. Hence it fails the de-authorization and non-substitution clauses of Definition A.3. \square

[Modal and Type-Theoretic Reading] One may equivalently regard $\text{Exec}(f)$ as a modal authorization $\Box f$ (“ f may be executed”) and refusal as the meta-act of forcing $\neg \Box f$. In a type-theoretic presentation, refusal has the shape

$$\rho : (\mathcal{W} \rightarrow \mathcal{R}) \rightarrow \perp,$$

A.4 A.2'' Minimality and Uniqueness of Refusal

[Object-Level Extension] An *object-level extension* of a governing mapping $f : \mathcal{W} \rightarrow \mathcal{R}$ is any mapping

$$f' : \mathcal{W} \rightarrow \mathcal{R}'$$

together with an injective embedding $\iota : \mathcal{R} \hookrightarrow \mathcal{R}'$ such that $f' = \iota \circ f$ on all previously admissible executions. Intuitively, f' augments the codomain with additional outputs, branches, labels, or confidence annotations while preserving executability.

[Minimality of Meta-Operational Refusal] Let $f : \mathcal{W} \rightarrow \mathcal{R}$ be a governing mapping. Among all mechanisms that can prevent the production of outputs in \mathcal{R} , refusal as defined in Definition A.3 is the unique mechanism that:

1. prevents execution without selecting an alternative output;
2. does not enlarge the representational space \mathcal{R} ;
3. preserves the abstraction status of f .

Proof. Consider any mechanism M that prevents certain executions of f .

Case 1: M operates by selecting an alternative output (e.g. a warning symbol, abstention token, or uncertainty estimate). Then M is an object-level extension of f in the sense of the preceding definition. Execution remains authorized; only the range of possible outputs is enlarged. Hence de-authorization does not occur, violating Definition A.3(1).

Case 2: M operates by conditionally switching to another rule $g : \mathcal{W} \rightarrow \mathcal{R}$. Then M is substitutional: it replaces one execution with another. Execution as such remains intact, violating Definition A.3(2).

Case 3: M operates by enlarging the internal state so as to encode contextual features required to decide whether to execute. Then M expands \mathcal{R} (or collapses abstraction entirely), violating Definition A.1 and destroying scalability.

The only remaining possibility is a mechanism that operates on the executability predicate $\text{Exec}(f)$ itself, without substituting outputs or enlarging \mathcal{R} . By Definition A.3, this is precisely meta-operational refusal. No weaker mechanism suffices, and no stronger mechanism preserves abstraction. \square

[Uniqueness] Refusal is the *unique minimal price* paid when abstraction retains autonomy: any attempt to preserve autonomy without meta-operational de-authorization either fails to halt execution or destroys abstraction.

[Why Uncertainty Is Not Refusal] Uncertainty quantification, confidence scoring, and probabilistic abstention are object-level extensions in the sense above. They produce additional representational content while leaving $\text{Exec}(f)$ intact. Consequently, they can modulate *what* is output, but not *whether* execution occurs. Uncertainty therefore cannot realize refusal.

[Human-in-the-Loop as External Refusal] A human-in-the-loop architecture does not internalize refusal into the abstraction. Rather, it externalizes ρ into a separate agent who alone possesses the capacity to negate $\text{Exec}(f)$. This confirms, rather than refutes, the externality of refusal.

[Modal Restatement] Let $\square f$ denote the authorization to execute f . Object-level safeguards implement transformations of the form

$$\square f \wedge \varphi \Rightarrow f'(w),$$

whereas refusal implements

$$\neg \square f,$$

with no consequent. The absence of a consequent is essential: refusal is not a transition but a suspension.

A.5 A.2'' Refusal as Non-Measurable Event

[Decision Space] Let $(\Omega, \mathcal{F}, \mu)$ be a measurable space of decision-relevant world states, and let

$$f : \Omega \rightarrow \mathcal{R}$$

be an abstraction executable whenever $\text{Exec}(f)$ holds.

[Measurable Safeguard] A safeguard mechanism is *measurable* iff it corresponds to a measurable function

$$s : \Omega \rightarrow \mathcal{R} \cup \mathcal{L},$$

where \mathcal{L} is a set of labels (e.g. warnings, abstentions, confidence scores).

[Non-Measurability of Refusal] Refusal, understood as de-authorization of execution, is not representable as a measurable event in $(\Omega, \mathcal{F}, \mu)$.

Proof. A measurable event $E \in \mathcal{F}$ corresponds to a subset of world states in which some output or label is produced. By Definition A.3, refusal does not produce an output in \mathcal{R} , nor a label in an extended codomain. Instead, it negates $\text{Exec}(f)$ itself.

Hence refusal does not correspond to a subset $E \subseteq \Omega$ on which f (or any extension of f) is evaluated. Rather, it removes the evaluation map entirely. Since measurability presupposes evaluation, refusal is not an event in \mathcal{F} .

Therefore refusal is non-measurable with respect to any probability measure μ defined over executable states. \square

[No Probability of Refusal] There exists no probability $\mu(\text{refusal})$ compatible with Definition A.3. Any attempt to assign such a probability necessarily reinterprets refusal as an object-level outcome.

[Why Risk Scores Cannot Encode Refusal] Risk estimation assigns probabilities to outcomes conditional on execution. Refusal negates the condition of execution itself. Consequently, no risk threshold, confidence interval, or uncertainty estimate can realize refusal without ceasing to be probabilistic.

A.6 A.2''' Undecidability of Refusal Detection

[Undecidability of Necessary Refusal] There exists no general algorithm that decides, for all abstractions f and all inputs $w \in \mathcal{W}$, whether refusal is required without violating the abstraction constraints of f .

Proof. Assume for contradiction that there exists an algorithm

$$D(f, w) \in \{\text{execute, refuse}\}$$

that correctly decides whether execution of $f(w)$ should be refused in all cases.

For D to decide refusal correctly, it must evaluate properties of the execution consequences of $f(w)$ that are not representable in \mathcal{R} , since refusal is triggered precisely by considerations external to the abstraction. Therefore D must determine non-local, non-represented properties of the computation induced by f .

This implies that D decides a semantic property of arbitrary programs not computable from their syntactic representation within the abstraction. By Rice’s Theorem, no such decision procedure exists.

Hence no general refusal-detection algorithm is possible without collapsing the abstraction back into the full world-state space \mathcal{W} . \square

[Impossibility of “Knowing When to Stop”] Any claim that a scalable abstraction can internally “know when to stop” is equivalent to claiming a solution to a non-trivial semantic decision problem over its own executions, and is therefore false.

[Speed as Structural Enemy of Refusal] Since refusal detection is undecidable, any real-time execution regime must approximate refusal by truncation, thresholds, or heuristics. These approximations preserve execution and therefore eliminate refusal by design.

[Relation to the Halting Problem] The present result is stronger than a halting-style argument. The halting problem concerns whether execution terminates. Refusal concerns whether execution should occur at all. The latter strictly subsumes the former and inherits its undecidability.

A.7 A.2'''' Synthesis

Taken together, Sections A.3–A.6 establish that refusal is:

1. meta-operational rather than executable,
2. non-substitutional rather than conditional,
3. non-measurable rather than probabilistic,
4. undecidable rather than computable.

These properties are not defects. They are the precise sense in which refusal constitutes the non-scalable core of autonomy.

A.8 Relation to Markov’s Principle

Markov’s principle (MP), in its computational formulation, asserts that for a decidable predicate P over \mathbb{N} ,

$$\neg\neg\exists n P(n) \Rightarrow \exists n P(n).$$

Intuitively, if it is impossible that a computation never succeeds, then it must succeed. While classically valid, this principle fails in general constructive settings precisely because it converts double-negated existence into positive execution.

The present framework reveals a strictly stronger obstruction. Refusal does not concern the existence of an output within a fixed domain, but the authorization of execution itself. Let $\text{Exec}(f)$ denote the executability of a governing mapping $f : \mathcal{W} \rightarrow \mathcal{R}$. An internalized “refusal-capable” abstraction would require an admissible inference of the form

$$\neg\neg \text{Exec}(f) \Rightarrow \text{Exec}(f),$$

thereby eliminating double negation at the level of executability.

However, by Sections A.3–A.6, executability is not a decidable predicate over world-states, nor is refusal an object-level event subject to enumeration, search, or probabilistic estimation. Refusal negates the evaluation map itself, rather than asserting the existence or non-existence of a value. Consequently, no Markov-style principle can be formulated internally for $\text{Exec}(f)$ without collapsing the abstraction.

In this sense, scalable abstractions fail to satisfy even the weakest admissibility conditions required to internalize Markov’s principle. The impossibility of endogenous refusal is therefore not merely a constructive limitation, but a structural consequence of abstraction itself. What Markov’s principle reveals at the boundary between constructive and classical logic, refusal reveals at the boundary between execution and governance.

A.9 Corollary: Extended Church’s Thesis and Executability

Extended Church’s Thesis (ECT) asserts that every total function is computable. In constructive arithmetic, ECT is often paired with Markov’s principle to recover classical exis-

tence results from double-negated claims.

The results of Sections A.3–A.8 show that neither assumption suffices to internalize refusal.

[Extended Church’s Thesis Does Not Recover Refusal] Assuming Extended Church’s Thesis, there still exists no internal procedure by which a scalable abstraction can realize refusal.

Proof. Extended Church’s Thesis concerns the computability of total functions on well-defined domains. Refusal, however, is not a function from \mathcal{W} to \mathcal{R} , nor a partial function, nor a total function with distinguished outputs. By Definition A.3, refusal negates the authorization of evaluation itself.

Therefore refusal is not a candidate for representation under ECT. Even if all total functions are computable, the executability predicate $\text{Exec}(f)$ is not a computable predicate internal to the abstraction. The obstruction identified here precedes questions of computability and persists under maximal computational assumptions. \square

This explains why appeals to “more powerful models,” “unbounded computation,” or “general intelligence” fail to address the refusal problem. The limitation is not computational but structural: refusal is not an uncomputed function but a non-functional suspension of execution.

A.10 Refusal and Bergsonian Duration

Henri Bergson distinguished *duration* (*durée*) from spatialized or metric time. Duration is not a sequence of discrete instants but a continuous, qualitative unfolding in which past, present, and future interpenetrate. Action within duration is not reducible to stepwise execution; it involves hesitation, reconsideration, and the retention of alternatives.

This distinction illuminates the formal results above. Scalable abstractions operate exclusively in spatialized time: execution is decomposed into discrete steps, states, or updates, each authorized by local rules. Such systems may pause, delay, or terminate, but these are transitions within execution, not suspensions of executability itself.

Refusal, by contrast, requires duration. To refuse is not merely to stop at a time-slice, but to withhold continuation in light of a temporally extended apprehension of context, consequence, and meaning. This apprehension cannot be discretized without loss. It is precisely what abstraction eliminates in order to achieve scalability.

[Speed Eliminates Refusal] Any system whose correctness depends on real-time or near-real-time execution necessarily eliminates refusal as a structural possibility.

Proof. By Section A.6, refusal detection is undecidable and cannot be resolved within bounded time. Real-time systems therefore approximate refusal through thresholds, heuristics, or truncation. These approximations preserve executability while sacrificing suspension. Since

duration cannot be compressed into bounded execution time without spatialization, refusal is excluded by design. \square

This clarifies why acceleration consistently correlates with ethical failure in automated systems. The problem is not insufficient foresight or misaligned values, but the elimination of duration itself. Where there is no duration, there can be no refusal.

A.11 Against the Dimensionalization of Refusal

A recurring error in the design and analysis of automated systems is the attempt to represent refusal as an internal state, variable, or coordinate within the system’s operational space. In such treatments, refusal is modeled as a distinct outcome, a threshold-crossing event, or a position along an axis of confidence or risk.

This maneuver is structurally analogous to treating time as a geometric dimension. In both cases, a condition of possibility is reified as an internal coordinate. The result is not clarification, but elimination. When refusal is dimensionalized—made enumerable, comparable, or interpolable—it ceases to function as refusal and becomes merely another mode of execution.

The defining feature of refusal is not its location within a state space, but its ability to suspend the application of that space altogether. Any representation that preserves the evaluative machinery while assigning refusal a position within it has already failed. Refusal cannot be indexed without being neutralized.

A.12 Temporal Suspension Versus Temporal Indexing

Execution unfolds in indexed time: a sequence of ordered steps, updates, or transitions. Such time is discrete, spatialized, and compatible with optimization. Suspension, by contrast, does not occur *at* a time-step, nor does it occupy an interval measured by the system clock. It interrupts the very progression by which steps are ordered.

Refusal is therefore not a temporal event but a temporal suspension. It cannot be represented as a delay, timeout, or pause without converting suspension into deferred execution. Delays preserve executability; refusal negates it.

This distinction explains why increasing temporal resolution—faster clocks, tighter feedback loops, real-time responsiveness—does not make systems safer. On the contrary, as execution becomes more finely indexed, the possibility of genuine suspension vanishes. What remains are only faster executions and more granular thresholds.

Refusal requires a form of time that is not reducible to indexation: a duration in which reconsideration, hesitation, and withdrawal are possible. Systems that operate exclusively in indexed time cannot host such duration and therefore cannot refuse.

A.13 Scale, Acceleration, and the Disappearance of Refusal

The loss of refusal is not an accidental byproduct of poor design, but a structural consequence of scale. As systems grow in scope and speed, they must replace suspension with approximation. Decisions must be made before their full context can be apprehended, and execution must proceed under conditions of partial knowledge.

At small scales, refusal appears viable because context remains legible and execution remains interruptible. As scale increases, interruption becomes costly, coordination becomes fragile, and hesitation is reframed as failure. What disappears is not judgment, but the capacity to withhold action without collapsing the system.

Acceleration intensifies this effect. Faster systems compress deliberation into execution and eliminate the temporal slack in which refusal could occur. This is why safeguards in high-speed domains take the form of thresholds, cutoffs, and automated overrides rather than suspensions. These mechanisms preserve flow while simulating caution.

The result is a paradox: the systems most in need of refusal are precisely those least capable of sustaining it. At scale, refusal must be externalized—or it will be eliminated entirely.

Synthesis

Refusal cannot be internalized without being destroyed. Attempts to represent it as a state, index it in time, or preserve it under acceleration all fail for the same reason: refusal is a suspension of execution, not a mode of it. Systems that scale by eliminating suspension eliminate refusal by necessity. Where refusal persists, it does so only by remaining external.

A.14 Proper Time, Coordinate Time, and the Error of Temporal Reification

It is important to distinguish the present argument from a common but misleading metaphor: that ethical failure in automated systems results from time “slowing down” or “speeding up.” No such dynamical claim is required, nor would it be correct.

In special relativity, time does not slow down in itself. What differs between reference frames is the decomposition of an invariant spacetime interval into temporal and spatial components. For two timelike-separated events, the proper time

$$\Delta\tau^2 = \Delta t^2 - \frac{\Delta x^2}{c^2}$$

is invariant. Coordinate time Δt varies with frame choice, and is minimized in the frame where the events occur at the same spatial location. Nothing physical happens to time; only the representational slicing changes.

This distinction clarifies the temporal structure of refusal. Refusal does not require that execution take longer, nor that time dilate. It requires that execution be suspended in a way that cannot be represented as a coordinate interval within the system’s operational frame. Delays, timeouts, and rate limits correspond to changes in coordinate time. They preserve executability. Refusal negates it.

Just as proper time is not an additional dimension but a geometric invariant that cannot be reduced to any single coordinate system, refusal is not an internal temporal state or duration that can be indexed within execution. It is a structural property of the system’s relation to action, not a measurable lapse between steps.

Attempts to internalize refusal as a timeout, pause state, or latency budget therefore commit the same category error as treating time itself as a manipulable dimension. In both cases, what is invariant is mistaken for what is representable. The result is not control but erasure: time becomes spatialized, and refusal becomes execution under another name.

A.15 Structural Agency

Definition A.3 (Structural Agency). A system exhibits structural agency if it satisfies all of the following:

1. (*Persistence*) Its effects propagate through time;
2. (*Generalization*) Its behavior applies across multiple contexts;
3. (*Constraint*) Its internal structure restricts downstream possibilities.

No requirement of intentionality, consciousness, or semantic understanding is imposed.

B Axis-Relativity Theorem

Theorem B.1 (Axis-Relativity of Superhuman Intelligence). For any cognitive artifact C , there exists at least one axis α along which C exceeds unaided human performance, and at least one axis β along which C is strictly inferior or blind.

Proof. Let C be defined by abstraction $A : \mathcal{W} \rightarrow \mathcal{R}$. By Definition A.1, A preserves only a subset of dimensions of \mathcal{W} . Along preserved dimensions (e.g., memory capacity, symbolic manipulation, numerical precision), performance scales independently of biological limits. Along discarded dimensions (e.g., perception, ethical judgment, contextual awareness), performance is identically zero.

Therefore superhuman capability along α necessarily implies subhuman incapacity along β . \square

Corollary B.2. “Artificial General Intelligence” understood as uniform superiority across all axes is a category error.

C Constitutive Limitation Theorem

Theorem C.1 (Constitutive Blindness of Abstraction). The limitations of an abstraction are necessary consequences of its scalability, not contingent engineering defects.

Proof. Assume an abstraction $A : \mathcal{W} \rightarrow \mathcal{R}$ satisfying Definition A.1. Reduction requires discarding degrees of freedom. Let $I(\mathcal{W})$ denote information content. Then

$$I(\mathcal{R}) < I(\mathcal{W}).$$

Recovering discarded information requires an inverse mapping A^{-1} , but such a mapping is undefined unless A is injective, contradicting reduction.

Thus blindness is structurally entailed by abstraction. \square

Corollary C.2. Any attempt to “add context” or “restore judgment” to an abstraction without sacrificing scalability is incoherent.

D Non-Intentional Agency Theorem

Theorem D.1 (Agency Without Intentionality). A system may exert causal agency without possessing beliefs, desires, or phenomenological awareness.

Proof. By Definition A.3, agency depends on persistence, generalization, and constraint. None of these properties require intentional states. Markets, legal codes, algorithms, and infrastructures satisfy all three conditions while lacking intentionality.

Therefore intentionality is not a necessary condition for agency. \square

Corollary D.2. Treating abstractions as “neutral tools” constitutes agency laundering.

E Opacity Necessity Theorem

Theorem E.1 (Necessity of Opacity). Total transparency of a scalable abstraction is mathematically impossible.

Proof. Let $A : \mathcal{W} \rightarrow \mathcal{R}$ be an abstraction. Transparency requires reconstructability of \mathcal{W} from \mathcal{R} . This requires A to be information-preserving.

But by Definition A.1, A is information-reducing. Therefore full transparency contradicts abstraction itself. \square

Corollary E.2. Opacity is an achievement condition of abstraction, not a political accident.

Corollary E.3. Local interpretability is compatible with abstraction; global interpretability is not.

F Externality of Responsibility Theorem

Theorem F.1 (Externality of Responsibility). Responsibility for the actions of an abstraction cannot be located within the abstraction itself.

Proof. By Definition A.2, refusal requires a meta-operation not derivable from the governing rule. By Theorem C.1, abstractions discard contextual degrees of freedom and therefore cannot generate meta-operations without ceasing to be abstractions.

Since responsibility requires the capacity to suspend or redirect execution in light of consequences, and abstractions lack refusal by design, responsibility cannot be internal to the abstraction.

Therefore responsibility must be external to the abstraction. \square

Corollary F.2. Attempts to assign moral responsibility to abstractions constitute category errors.

Corollary F.3. Responsibility must be imposed through governance, law, or human intervention external to the abstraction.

G Invariant Alignment Theorem

Theorem G.1 (Invariant Alignment Theorem). Alignment of scalable abstractions must operate at the level of representational invariants rather than behavioral control.

Proof. Behavioral alignment presupposes that a system can conditionally refuse or redirect execution based on context. By Theorem C.1 and Definition A.2, contextual refusal is incompatible with scalable abstraction.

Let \mathcal{R} be the representational space of abstraction A . Alignment is achievable iff the set of misaligned states $M \subset \mathcal{R}$ is unreachable under all permissible transformations of A .

Therefore alignment must consist in constraining the topology of \mathcal{R} such that M is excluded by construction. \square

Corollary G.2. Behavioral safeguards without representational constraints are necessarily brittle.

Corollary G.3. Invariant alignment is a mathematical necessity, not an engineering preference.

H Extinction Impossibility Theorem

Theorem H.1 (Impossibility of the Detached Maximizer). A detached, refusal-capable, adversarial maximizer is ontologically incompatible with scalable abstraction.

Proof. A detached maximizer must possess:

1. autonomous goal revision,

2. contextual refusal,
3. global situational awareness.

Each property requires retention of degrees of freedom explicitly discarded by abstraction (Definition A.1). Therefore any system possessing these properties is not scalable abstraction.

Extinction narratives that posit such agents presuppose a contradictory ontology. \square

Corollary H.2. Existential risk does not arise from adversarial intent but from inexorable execution within ecological and institutional flows.

Corollary H.3. Containment, not persuasion, is the appropriate safety paradigm.

I Governance Theorem

Theorem I.1 (Governance as Exogenous Refusal). Effective governance of abstractions consists in the deliberate reintroduction of refusal external to the abstraction.

Proof. By Theorem F.1, responsibility cannot be internal. By Theorem G.1, alignment must constrain representability rather than behavior.

Therefore governance must introduce:

1. veto mechanisms,
2. shutdown capabilities,
3. institutional override authority,
4. enforced boundary conditions.

Each constitutes refusal imposed from outside the abstraction. \square

Corollary I.2. Kill switches and veto power are not failures of design but necessary complements to abstraction.

Corollary I.3. Institutions that eliminate refusal eliminate their own capacity for self-correction.

J Meta-Theorem: Closure of the System

Theorem J.1 (Closure and Consistency of the Refusal Framework). The preceding theorems form a closed, non-contradictory system describing abstraction, agency, alignment, and governance.

Proof Sketch.

1. Abstraction requires reduction (Theorem C.1).

2. Reduction amputates refusal (Definition A.2).
3. Amputation produces non-intentional agency (Theorem D.1).
4. Agency without refusal necessitates external responsibility (Theorem F.1).
5. External responsibility requires invariant alignment (Theorem G.1).
6. Failure to impose invariants yields ecological risk, not adversarial takeover (Theorem H.1).
7. Governance restores refusal exogenously (Theorem I.1).

No theorem contradicts another; each depends on the previous. The system is therefore internally consistent and complete with respect to its domain. \square

Corollary J.2. Any proposal that demands contextual judgment, moral agency, or endogenous refusal from scalable abstractions violates at least one theorem in this system.

Corollary J.3. The autonomy of refusal is the non-scalable remainder upon which all ethical governance depends.

References

- [1] H. Melville, *Bartleby, the Scrivener*, 1853.
- [2] H. D. Thoreau, *Civil Disobedience*, 1849.
- [3] G. Agamben, *The Use of Bodies*, Stanford University Press, 2015.
- [4] H. Martinson, *Aniara*, 1956.