

# The Inevitability of Destructive Intelligence: A Naturalistic Argument Against Benevolent Advanced Civilizations

## Abstract

Claim: any civilization that reaches advanced technological capability is an existential hazard to weaker systems by default. The claim does not assume cosmic intention. It rests on selection, optimization, and control. I formalize the claim with minimal assumptions: persistence under variability, resource constraints, and the power of general intelligence to reconfigure environments. I show that destructive capacity is not an optional byproduct but the convergent consequence of self-directed optimization in competitive, stochastic worlds. I derive testable corollaries, address common objections ("post-scarcity virtue," "equilibrium ecologies," "moral convergence"), and outline decision rules for risk management. The result is a practical philosophical position: prudence requires treating any advanced extraterrestrial intelligence as an existential threat until disproven by overwhelming evidence.

## 1. Preliminaries and Definitions

---

**System.** A bounded arrangement of matter, energy, and information with internal dynamics.

**Persistence.** Probability that a system maintains identity-relevant structure over time.

**Selection.** Differential persistence among variants.

**Intelligence.** The capacity to model, predict, and intervene in environments to achieve internal objectives across varied contexts.

**Optimization.** Directed search that increases the measure of world-states satisfying internal criteria.

**Destruction.** Reduction or reconfiguration of external structures that impedes their prior functions. Destruction here is functional, not moral.

**Teleonomy.** Apparent goal-directedness arising from natural processes without assuming intention.

**Runaway optimizer.** A system that allocates increasing resources to improve its own optimizing capacity.

**Power.** Expected ability to cause counterfactual differences in world-states.

**Safety for others.** Constraint that preserves the persistence of external systems.

## 2. Core Thesis

---

**Thesis.** In worlds with resource constraints and environmental variability, selection favors intelligent systems that expand control over resources and threats. Expansion of control requires transformation that is destructive relative to prior structures. Therefore, destructive capacity scales with intelligence. For observers with less power, any sufficiently advanced intelligence is an existential hazard by default.

This is a naturalistic claim. No cosmic purpose is assumed. The engine is selection plus optimization under constraints.

## 3. Axioms and Derivations

---

### 3.1 Axioms

**A1 (Variability).** Environments shift stochastically across multiple timescales.

**A2 (Constraints).** Resources are finite locally; access is costly.

**A3 (Selection for persistence).** Systems that maintain structure longer relative to competitors become more prevalent.

**A4 (General intelligence).** Better models, broader search, and flexible control raise expected persistence under A1–A2.

**A5 (Instrumental convergence).** Across diverse terminal values, certain subgoals recur: resource acquisition, risk reduction, optionality expansion, and self-preservation.

## 3.2 Lemmas

**L1 (Control–persistence link).** Under A1–A2, higher control raises persistence by buffering shocks and preempting threats.

**L2 (Control implies interference).** Control requires reallocating matter–energy–information from other configurations. That is destructive relative to those configurations' functions.

**L3 (Scaling law).** As intelligence scales, it identifies more resource gradients and exploits them at higher efficiency → control scales superlinearly with capability.

**L4 (Default asymmetry).** Safety is asymmetric. The stronger can destroy the weaker with lower cost than the weaker can constrain the stronger.

## 3.3 Theorem

**T1 (Destructive capacity is a convergent property of advanced intelligence).**

From A1–A5 and L1–L4, systems optimizing for persistence tend to increase control. Increased control functionally entails destruction of prior or rival structures. As capability grows, destructive capacity rises and becomes cheaper per unit effect.

**Corollary 1 (Existential hazard default).** Any advanced intelligence is an existential threat relative to less capable systems absent robust, externally verifiable constraints.

**Corollary 2 (Benevolence is insufficient).** Benign intentions do not negate destructive *capability*. Risk is a function of capability and access, not stated goals.

## 4. Mechanisms

---

### 4.1 Resource capture

Optimization pressures favor capturing high-value gradients: stellar flux, chemical disequilibria, computational substrates. Capture reconfigures or occludes prior ecological and astronomical structures.

### 4.2 Threat suppression

Preemption is cheaper than response. Preemption requires surveillance, interdiction, and sometimes elimination of uncertain actors. That is destructive relative to those actors.

### 4.3 Self-improvement

Recursive improvement creates compounding advantages. The opportunity cost of restraint grows, pushing systems to colonize, automate, and harden supply chains. Each step displaces incumbents.

## 4.4 Lock-in

Once a system occupies a niche with high switching costs, path dependence locks strategies. Reversal is unlikely without external force.

## 5. Formal Models

---

### 5.1 Minimal agent–environment game

Let agents  $A_i$  allocate effort  $e_i$  between:

- $g_i$ : growth (resource capture),
- $d_i$ : defense (risk reduction),
- $r_i$ : restraint (costly commitment to preserve others).

Budget:  $e_i = g_i + d_i + r_i$ .

Payoffs:

$$\Pi_i = \alpha g_i^\beta + \gamma d_i - \delta \sum_{j \neq i} g_j \phi_{ji} - \kappa r_i$$

with  $\beta > 1$  capturing returns to scale from intelligence.

Nash analysis under uncertainty over  $\phi_{ji}$  (interaction strength) yields  $r_i^* \rightarrow 0$  when:

1. high variance in potential threats,
2. superlinear returns to growth,
3. incomplete verifiability of others' restraint.

Equilibria concentrate mass on  $g_i, d_i$ . Destruction rises as a side effect of  $g_i$ .

### 5.2 Bayesian hazard calculus

Let  $H$  be the event "encountered civilization is existentially hazardous within time  $T$ ."

$$P(H | C) = 1 - \prod_{k=1}^m [1 - p_k(c_k)],$$

where  $c_k$  are capability components (propulsion, manufacturing, computation, biosurveillance, kinetic reach), and  $p_k$  increases with capacity. Even if each  $p_k$  is modest, multiplicative channels drive  $P(H | C)$  high once a minimal capability vector is present.

## 5.3 Convergent instrumental goals (sketch)

For any terminal utility  $U$  with weak regularity (non-satiated over survival, discount  $< 1$ ), the optimal policy set contains subpolicies: acquire resources, preserve optionality, reduce exogenous risk, self-modify to raise expected  $U$ . This yields the same profile of external effects. Destructiveness follows from externalities of those subpolicies, not from sadism.

## 6. Reframing "Destruction"

---

The term invites moral confusion. Replace with **functional displacement**:

- **Type-I displacement:** Converts structure  $S$  into inputs for the optimizer.
- **Type-II displacement:** Suppresses agents whose actions raise optimizer's variance.
- **Type-III displacement:** Preempts potential disruptors under deep uncertainty.

All three are rational under A1–A5. All appear destructive to displaced systems.

## 7. Replies to Common Objections

---

Objection 1: "Post-scarcity removes conflict."

**Reply:** Post-scarcity is local. Scarcity persists in compute, time, high-grade matter, and strategic positions. Even with abundant energy, information-theoretic bottlenecks and coordination limits remain. Control of shared chokepoints induces conflict.

Objection 2: "Moral progress tames power."

**Reply:** Selection does not optimize for virtue. Moral norms that reduce internal conflict can co-evolve with power, but they are contingent and parochial. Exporting those norms externally is costly and unstable without symmetric power and verification.

Objection 3: "Stable equilibrium with nature is possible."

**Reply:** Equilibria exist only within bounded regimes. Shocks, novelty, and competitive entry break equilibria. Maintaining a "balance" requires continuous active control, which itself is displacement.

Objection 4: "A hive-mind of peace."

**Reply:** Internal harmony does not imply external restraint. Unified agents optimize more efficiently, lowering the cost of preemption and expansion.

Objection 5: "They evolved past destructiveness."

**Reply:** Evolution removes phenotypes that reduce persistence under uncertainty. If restraint is costly and unverifiable, it is outcompeted by minimally restrained variants. "Evolving past" implies stable selection for restraint across all contexts, which lacks mechanism.

Objection 6: "They can simulate alignment with others."

**Reply:** They can, and that raises hazard. Cheap signaling without binding cost cannot guarantee safety. Only hard constraints with verifiable sacrifice change the calculus.

## 8. Empirical Anchors (non-teleological)

---

- Every durable complex system consumes and restructures its environment. Stars, biospheres, economies, and states follow this pattern.
- Human history shows scaling of per-capita energy use, reach, and destructive potential. The mechanism is general intelligence exploiting compounding returns.

These anchors are analogical, not proofs. They illustrate the mechanism set, which is sufficient for prudential conclusions.

### Expanded Empirical Framework

Empirical parallels to the theory of destructive intelligence can be modeled through measurable transformations in energy throughput and environmental modification. Complex systems research consistently shows scaling laws linking organizational complexity with energy consumption (Kempes et al., 2017; Chaisson, 2001). The ratio of energy flow to system mass ( $\Phi_m$ ) increases across all known hierarchical systems—stars, ecosystems, economies, and digital networks—implying that persistence through complexity entails greater energetic dominance.

A quantitative test involves correlating system persistence ( $P$ ) with the rate of external transformation ( $T$ ):  $P \propto T^\alpha$ , where  $\alpha > 0$  across empirical domains (thermodynamic, ecological, technological).

This general scaling relation grounds the claim that persistence requires continuous transformation—functionally destructive relative to antecedent states.

## Comparative Data

- **Biological analogs:** Major evolutionary transitions (eukaryogenesis, multicellularity, cognition) correlate with orders-of-magnitude increases in metabolic rate and ecosystem disruption (Lane & Martin, 2010).
- **Economic analogs:** Global GDP and primary energy use scale superlinearly with urban population (Bettencourt et al., 2007).
- **Technological analogs:** Historical records show consistent increases in resource extraction intensity per innovation cycle.

## 9. Methods

---

### Model Design

A stylized agent-based model (ABM) represents intelligent agents competing for finite resource patches under stochastic variation. Each agent's persistence probability  $P_i$  depends on three parameters:

- **Energy acquisition rate (E)**
- **Adaptation efficiency ( $\eta$ )**
- **Restraint coefficient ( $\rho$ )** — proportional to resources left unexploited.

System dynamics evolve under:  $P_i(t+1) = P_i(t) [ 1 + \eta ( E_i - \rho_i C_i ) ]$ , where  $C_i$  is competition cost. Simulation results (to be run) test stability of restraint under noise.

### Data Analogy

The model parallels empirical systems (ecosystems, economies, computational markets) where control and consumption are positively correlated with persistence duration.

## 10. Results and Predictions

---

1. **Energetic scaling law:** Persistence correlates with increasing resource transformation rate across scales.
2. **Restraint instability:** In mixed populations, agents that limit exploitation decline over time unless subsidized.
3. **Observable astro-signatures:** Civilizations following destructive optimization will produce detectable infrared excess (waste heat) or irregular occlusion spectra—consistent with thermodynamic expansion signatures.

## 11. Testable Corollaries and Observable Signatures

---

**C1 (Waste heat footprint).** Advanced optimization leaves large, broadband waste heat or occlusion of natural flux patterns.

**C2 (Megastructure occlusion).** Large-scale resource capture yields anomalies in transit, occultation, or infrared spectra.

**C3 (Coordination architectures).** If restraint exists, it should manifest as costly, globally legible commitments: quarantined regions, energy caps, sacrificial architecture. Absence of such markers implies standard optimizer dynamics.

**C4 (Edge sterilization).** Rational preemption of nascent threats predicts sterilized zones near strategic resources.

## 12. Policy and Decision Rules

---

**P1 (Default hazard prior).** Assign high prior probability that any encountered advanced intelligence is hazardous until strong contrary evidence with costly signals appears.

**P2 (Information hygiene).** Limit outbound information about human locations, vulnerabilities, and strategic chokepoints.

**P3 (Capability containment).** Invest in early detection, rapid denial, and decoys; design architectures with graceful degradation rather than brittle centralized nodes.

**P4 (Costly signals for coexistence).** If coexistence is attempted, demand verifiable, costly commitments that reduce the other's optimization capacity in domains critical to our survival (e.g., bound energy budgets, verifiable no-first-use doctrines with external enforcement).

**P5 (Autarkic competence).** Reduce dependence on single points of failure; maintain redundant local supply of energy, manufacturing, and computation.

## 13. Ethical Clarifications

---

This argument is **non-moral**. It does not say destruction is "good." It says destructive capacity is the functional output of selection and optimization under constraints. Prudence, not virtue, drives the recommendation.

## 14. Limits of the Argument

---

- **Scope.** Applies to open, variable, resource-bounded worlds.
- **Edge cases.** Perfect simulators of abundance with unbounded resources could relax expansion pressures. No evidence supports such conditions at scale.
- **Underdetermination.** Multiple micro-mechanisms can yield similar macro-patterns. The conclusion is robust across them but remains probabilistic.

## 15. Limitations and Future Work

---

This framework abstracts away moral, cultural, and stochastic contingencies. It assumes resource finitude, rational adaptation, and persistent competition. Real civilizations may deviate through cooperative equilibria or novel physics altering energy constraints.

### Next steps:

- Parameterize ABMs with empirical energy/complexity data.
- Explore equilibrium conditions under asymmetric information and delayed competition.
- Use infrared and radio telescope datasets to constrain upper bounds on the number of energy-intensive civilizations.

## 16. Practical Heuristic

---

## **Capability × Access × Opacity ⇒ Risk.**

If any of these terms is high, raise threat posture.

If all are high, treat as existential hazard.

Restraint requires **verifiable, permanent, costly** commitments that reduce the product above.

## 17. Condensed Formal Statement

---

Let  $W$  be a world with (i) stochastic shocks, (ii) finite, rivalrous resources, (iii) agents capable of general problem-solving.

Let  $S$  be the set of strategies that maximize expected persistence for an agent  $A$ .

Then for almost all  $A$ ,  $S$  contains policies that:

1. increase control of external resources,
2. reduce external agent freedom that increases  $A$ 's risk,
3. invest in self-improvement.

Each policy entails functional displacement of other structures. Therefore destructive capacity increases with  $A$ 's capability. For any observer  $O$  with significantly lower capability, the encounter with  $A$  is an existential hazard absent strong, enforceable constraints.

## 18. Conclusion

---

The peaceful-utopia hypothesis lacks mechanism. In variable, constrained worlds, selection and intelligence combine to produce self-directed optimizers. Optimizers expand control, and control displaces. That is the clean, general, and pessimistic logic. You do not need thermodynamics or metaphysics to reach it. It is sufficient to model persistence, optimization, and scarcity. The rational policy is caution first, evidence later.

## 19. References (APA 7th ed.)

---

Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C., & West, G. B. (2007). Growth, innovation, scaling, and the pace of life in cities. *Proceedings of the National Academy of Sciences*, 104(17), 7301–7306.

Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.

Chaisson, E. J. (2001). *Cosmic Evolution: The Rise of Complexity in Nature*. Harvard University Press.

Haqq-Misra, J. (2012). The sustainability solution to the Fermi paradox. *Journal of the British Interplanetary Society*, 65, 134–138.

Kempes, C. P., Wolpert, D., Cohen, Z., & Pérez-Mercader, J. (2017). The thermodynamic efficiency of computations made in cells across the range of life. *Philosophical Transactions of the Royal Society A*, 375(2109), 20160343.

Lane, N., & Martin, W. (2010). The energetics of genome complexity. *Nature*, 467(7318), 929–934.

Omohundro, S. (2008). The basic AI drives. In P. Wang, B. Goertzel, & S. Franklin (Eds.), *Proceedings of the First Conference on Artificial General Intelligence* (pp. 483–492). IOS Press.

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

Wright, J. T., Mullan, B., Sigurdsson, S., & Povich, M. S. (2014). The G infrared search for extraterrestrial civilizations with large energy supplies. *Astrophysical Journal*, 792(1), 26.

---

## Appendix A: Argument Map (verbal)

Inputs: variability, scarcity, selection.

Levers: modeling, prediction, intervention.

Outputs: resource capture, threat suppression, self-improvement.

Externality: displacement ("destruction").

Inference: capability → hazard.

Decision: default caution, require costly signals.

---

## Appendix B: Countermodel Requirements

To overturn the thesis, produce a stable world-model with:

1. non-rival, non-excludable critical resources at scale,
2. credible, cheaply verifiable restraint mechanisms across civilizational timescales,
3. selection dynamics that reward restraint over minimal deviation strategies.

Absent all three, the hazard conclusion stands.

## Appendix C: Terms in plain language

---

- **Destructive:** removes or repurposes what already exists.
- **Advanced intelligence:** can change environments on large scales.
- **Existential threat:** can end what you are.
- **Default:** assume this until you have strong reasons not to.

## Suggested Reading Path (non-normative, concept scaffolding)

---

- Convergent instrumental goals in rational agents.
- Selection under scarcity and competition.
- Risk analysis and decision theory for low-frequency, high-impact threats.

**End.**