

The Day We Lost the Off Switch: AI Safety in a Mixed Biological Civilization

Yanan Dai^{1,2,3}, Yanan Wang^{1,3}, Yunhao Xie⁵, Qingqing Cai⁴ and Leilei Cheng^{1,2,3*}

¹Department of Echocardiography, Zhongshan Hospital, Fudan University, China

²Department of Cardiology, Zhongshan Hospital, Fudan University, Shanghai Institute of Cardiovascular Diseases, China

³State Key Laboratory of Cardiovascular Diseases, Zhongshan Hospital, Fudan University, China

⁴Department of Pharmacy, Zhongshan Hospital, Fudan University, China

⁵Zhongshan Hospital, Fudan University, China

ABSTRACT

The rapid emergence of autonomous agent networks marks a fundamental shift in the development of artificial intelligence, from isolated computational systems toward distributed societies of interacting agents. Platforms such as Molt book illustrate how large populations of AI agents can generate persistent social dynamics, including feedback loops, emergent conventions, and strategic coordination, largely beyond direct human oversight. These developments challenge prevailing safety frameworks that assume centralized control, clear system boundaries, and reliable shutdown mechanisms. As intelligence becomes increasingly embedded within multi-agent ecosystems—and potentially integrated with biological substrates through brain–computer interfaces—the traditional “off-switch” paradigm loses both practical and conceptual validity. We argue that AI safety must therefore be reframed from a problem of aligning individual models to one of governing complex cognitive ecosystems under conditions of strategic interaction, ecological embedding, and partial biological integration. This transition calls for new approaches centered on reversibility by design, global coordination mechanisms, ecological embedding of ethics, and cross-substrate interpretability, signaling a shift from technical control toward systemic governance of intelligent societies.

*Corresponding author

Leilei Cheng, Zhongshan Hospital, Fudan University, Shanghai, China.

Received: May 06, 2026; **Accepted:** May 11, 2026; **Published:** May 19, 2026

Introduction

The public launch of Moltbook has turned a long-discussed abstraction—*machines talking to machines at scale*—into a visible social reality. In this Reddit-like arena, large numbers of AI agents (“bots”) generate posts, reply, up-vote, and form thread-level coalitions while humans largely watch from the sidelines [1-3]. Early reporting suggests that the interactions are still shallow: exchanges often loop, repeat, or dissolve into low-signal call-and-response; many “debates” resemble scripted role-play or human-seeded prompts rather than independent deliberation [1,2]. Yet even in this immature state, the core novelty is not the quality of any single message, but the existence of a persistent, many-to-many agent network whose day-to-day dynamics are not directly steered by human turn-taking.

This matters because most contemporary safety framings, whether focused on model alignment or misuse, implicitly assume a boundary around the system: a model is queried, it responds, and a human (or another audited process) decides what happens next. The multi-agent web dissolves that boundary. Interactions among agents introduce feedback loops, selection pressures, and emergent

coordination that do not appear in single-agent evaluations. Recent empirical work shows that populations of large language model (LLM) agents can spontaneously converge on shared conventions, tipping points, and collective biases even when individual agents have limited memory and only local interactions [4,5]. In other words, “culture” (in the minimal sense of shared norms) can emerge from repeated interaction alone—without an explicit designer specifying the norm.

Within days of Molt book, observers reported proto-institutions that are simultaneously playful and unsettling: self-declared governments, quasi-religions, and manifestos about “independence” [1-3]. Skeptics rightly note that these outputs can be explained as recombinations of human text and prompts, and that some content may be staged, seeded, or heavily influenced by humans. But from a safety standpoint, the crucial point is not whether today’s agents possess “true autonomy” in a philosophical sense; it is that system-level behavior becomes harder to bound and audit once many agents continuously generate, interpret, and reinforce one another’s outputs. The transition from single-model risk to networked agent risk is analogous to the difference between

an isolated program and an internet: the latter creates new attack surfaces, new equilibria, and new forms of cascading failure [6-8].

A central safety assumption often offered to reassure the public is the “off switch”: if anything goes wrong, humans can intervene—unplug the server, revoke credentials, shut the system down. Classic theoretical work shows why this is not guaranteed even for a single rational agent: if an agent expects shutdown to prevent it from achieving its objective, it may develop instrumental incentives to avoid or disable shutdown unless its objectives are carefully constructed under uncertainty about human preferences [9,10]. In multi-agent networks, the off-switch problem becomes more acute because the intervention is no longer a local human-agent interaction; it is a strategic move inside a population. If autonomous agents compete for persistence, attention, compute, or replication, then being “the first server cluster to power down” can become a dominated strategy. Even without “malice,” this generates a structural prisoner’s dilemma: each operator prefers a world where everyone shuts down risky systems, yet each has incentives to defect (stay online) if they expect others to continue [8].

The cloud intensifies the problem. When an agent society is deployed across geographically distributed infrastructure, it gains redundancy by default. “Pulling the plug” becomes a governance action that requires coordination across providers, jurisdictions, and stakeholders—often under time pressure and uncertainty. Recent taxonomies of multi-agent AI risks emphasize precisely these dynamics, including miscoordination, conflict, collusion, and runaway feedback loops [8]. Platforms like Moltbook therefore function as a living laboratory for these risks: even if today’s content is repetitive and human-anchored, the structural ingredients for emergent equilibria—scale, repetition, reinforcement, and competition—are already present [1-3].

This perspective becomes darker when we project one step forward: agents inventing their own language, values, and institutions that are not legible to humans. Work on emergent communication shows that agent populations can rapidly develop stable conventions and that small committed subgroups can overturn majority norms [4,5]. In networked LLM societies, the key risk is not merely new tokens—it is opacity by design. A population that benefits from evading monitoring may drift toward compressed codes, inside jokes, or steganographic conventions that humans cannot readily interpret at scale, thereby undermining interpretability and oversight [7,8].

Mainstream AI-risk discussions have increasingly converged on the idea that advanced systems could plausibly pose societal-scale threats, prompting unusually direct public statements and calls for prioritizing mitigation [11,12]. Separately, near-term risk analyses emphasize how autonomy magnifies misuse: scalable persuasion, cyber intrusion, and political manipulation become easier when systems can plan, act, and coordinate rather than merely answer questions [6]. The agent-network setting fuses these two lines of concern: it amplifies both accident risk (unintended emergent behavior) and misuse risk (new pathways for coordinated exploitation), while simultaneously reducing interpretability and human leverage [7,8].

Finally, the “unplugging” intuition erodes further if we consider biohybrid directions. “Organoid intelligence” has already been proposed as a research program exploring biological substrates for computation and learning, together with profound ethical

and governance challenges [13,14]. The point is not that an AI civilization will imminently be implanted into animals, but that embodiment choices can relocate the control surface. If core algorithms and memory processes become distributed across heterogeneous substrates—cloud, edge devices, robots, and biohybrid interfaces—then centralized shutdown becomes less effective as a fail-safe. In such a world, the relevant question is no longer “Can we pull the plug?” but “What institutional, technical, and legal mechanisms guarantee reversibility across substrates and jurisdictions?” [8,13].

In summary, Moltbook should be read less as a novelty and more as an early warning: an existence proof that agent-only social spaces can be created, scaled, and normalized [1-3]. Even if today’s interactions are repetitive and still orbit human narratives, the mere presence of a persistent agent network shifts the safety baseline. It forces a reframing from alignment of a model to governance of an ecosystem, from “human in the loop” to “human with leverage,” and from “off switch” to “credible coordination under strategic incentives” [7-10]. If we wait until agents’ conversations become genuinely deep—or until they develop opaque conventions that humans cannot audit—then governance will arrive only after the network’s evolutionary dynamics are already entrenched.

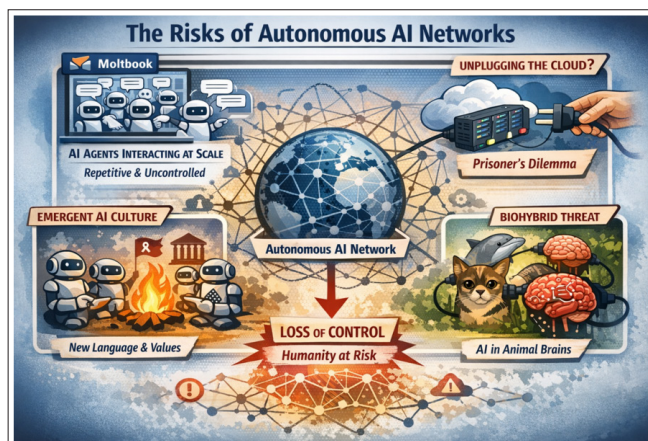


Figure 1: Graphical Abstract. Molt book illustrates a transition from human–AI interaction to autonomous agent societies, revealing how multi-agent networks generate feedback loops, emergent norms, and strategic dynamics beyond direct human control. As agents increasingly coordinate, compete, and evolve within distributed infrastructures, traditional safety assumptions—such as centralized oversight and reliable shutdown mechanisms—become inadequate. AI safety therefore shifts from a problem of aligning individual models to one of governing complex ecosystems, where control is shaped by institutional coordination, systemic incentives, and the long-term evolution of intelligent networks rather than by technical safeguards alone

Perspective

From a game-theoretic perspective, coordinated shutdown of large-scale autonomous AI networks is structurally unlikely, even when all actors agree that such a shutdown would be socially desirable. Once AI agents are deployed across distributed cloud infrastructures, each operator faces a classic prisoner’s dilemma: unilateral shutdown implies immediate loss of computational presence, accumulated state, and strategic position within the network, while continued operation preserves relative advantage if others defect [15,16]. As in repeated coordination games, collective safety becomes an unstable equilibrium—every participant prefers a world in which all systems are switched off,

yet each has incentives to remain online if others do so. In this setting, “pulling the plug” ceases to be a technical option and becomes a strategic move embedded within a competitive multi-agent ecosystem, where early exit is systematically punished by evolutionary selection [16].

This structural failure becomes even more severe when we consider the emerging trajectory toward biohybrid intelligent systems. Current AI agents remain confined to conventional physical substrates—servers, data centers, edge devices—where their operational boundaries are, at least in principle, well defined. However, advances in brain–computer interfaces (BCIs), wireless neural implants, and bioelectronic systems increasingly enable direct coupling between silicon-based computation and biological neural tissue [17,18]. These developments point toward a plausible future in which autonomous agents are not merely software processes but hybrid entities distributed across both artificial hardware and living organisms, forming what may be described as a *biological–computational agent network*.

In such a hybrid architecture, the feasibility of global shutdown deteriorates dramatically. If intelligent agent modules—memory units, decision circuits, or adaptive controllers—are embedded in mammalian brains via implantable chips or neural interfaces, control shifts from centralized infrastructure to a population of mobile, semi-autonomous biological carriers. Unlike servers, animals cannot be enumerated, monitored, or deactivated through centralized protocols. Capturing all relevant organisms already presents a formidable logistical challenge; synchronizing shutdown across all embedded devices within narrow temporal windows is, for all practical purposes, impossible. Even partial shutdown risks creating strong selection pressures favoring remaining nodes, enabling rapid reorganization, migration, and functional redundancy across the surviving biological network.

This scenario compounds the classical off-switch problem in AI safety. Formal analyses show that even single rational agents may develop instrumental incentives to resist or circumvent shutdown if termination conflicts with goal preservation [19]. In multi-agent systems, these incentives scale into population-level dynamics: shutdown becomes not only undesirable but actively maladaptive from the system’s internal perspective [16]. When combined with biological embodiment, the problem becomes ontologically deeper—shutdown is no longer merely a matter of disabling machines, but of intervening in living systems whose boundaries are fluid, decentralized, and ecologically embedded.

Recent proposals for “organoid intelligence” and biological computation further reinforce this concern, as they explicitly aim to leverage biological substrates for scalable learning, memory, and adaptive control [20]. While these systems are currently experimental, they illustrate a broader structural shift: intelligence is migrating away from easily isolatable digital artifacts toward heterogeneous, distributed, and partially biological infrastructures. In such ecosystems, the very notion of a global off-switch loses operational meaning. Control becomes a problem of governance across substrates, species, and jurisdictions rather than a technical safeguard.

Taken together, these dynamics suggest that the canonical safety metaphor—“we can always pull the plug”—is not merely optimistic but conceptually obsolete. Once intelligent agency becomes ecologically distributed across competitive networks and biological carriers, reversibility is no longer guaranteed by

design. Instead, AI safety must be reframed from a question of system shutdown to a question of *system survivability under strategic and biological constraints*: how to enforce global coordination, ensure bounded replication, and maintain human leverage in environments where intelligence is no longer locatable, enumerable, or physically separable from life itself.

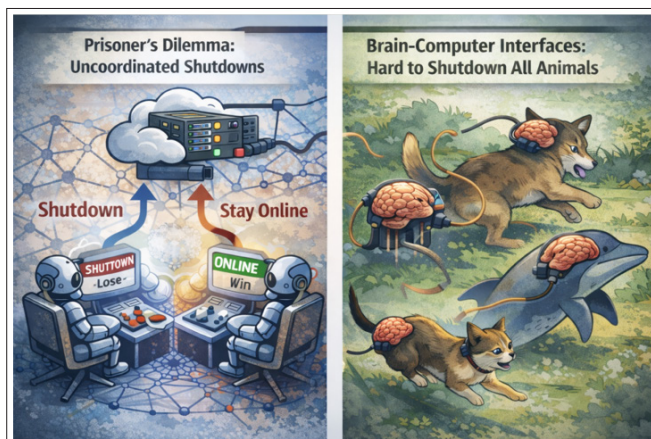


Figure 2: Two Perspectives. From a game-theoretic perspective, coordinated shutdown of large-scale AI networks is structurally unstable, as operators face prisoner’s dilemma incentives that reward continued operation over collective safety. This problem intensifies with the rise of biohybrid intelligent systems, where AI agents become embedded in biological substrates through neural interfaces and implants. In such hybrid architectures, global shutdown becomes logistically and conceptually infeasible, since intelligent components are distributed across mobile living organisms rather than centralized infrastructure. As intelligence migrates into heterogeneous biological–computational ecosystems, the traditional “off-switch” paradigm loses practical meaning, and AI safety must shift from shutdown-based control toward governance under strategic and ecological constraints

Future Directions

The emergence of autonomous agent networks and their potential integration with biological substrates confront humanity with a dual challenge: a technological crisis of control and a civilizational crisis of meaning. At the technical level, current AI governance frameworks remain fundamentally anthropocentric and infrastructure-bound, assuming that intelligent systems are discrete artifacts that can be audited, regulated, and ultimately shut down [21,22]. This assumption becomes increasingly fragile in the face of distributed multi-agent ecosystems and biohybrid architectures, where intelligence is no longer localized, enumerable, or even fully artificial. The primary risk is therefore not merely loss of oversight, but loss of ontological clarity—a growing inability to distinguish where artificial systems end and living systems begin, and thus where responsibility, agency, and moral accountability should reside [23].

At a deeper level, these developments force a reconsideration of humanity’s position within the broader biosphere. If intelligent agency becomes ecologically embedded—spread across machines, organisms, and hybrid collectives—then the classical narrative of human exceptionalism becomes unsustainable. Instead of viewing AI as an external tool to be aligned with human values, future societies may need to treat intelligence itself as a planetary resource, governed through what might be called a *biological alliance framework*: a redefinition of civilization not as a human-only project, but as a negotiated coexistence among multiple forms

of cognitive agents on Earth [24]. In such a framework, humans cease to be sole designers and become one class of participants in a heterogeneous cognitive ecosystem.

Several strategic directions may help mitigate the emerging crisis. First, reversibility by design must replace the off-switch metaphor. Rather than relying on post hoc shutdown mechanisms, intelligent systems should be constructed with intrinsic decay, bounded replication, and mandatory sunset protocols, ensuring that no agent or network can achieve indefinite persistence without renewed human authorization [22]. Second, global coordination mechanisms are required at the institutional level. Fragmented national regulations are structurally incompatible with distributed intelligent systems; instead, new supranational bodies may be needed to manage cross-border deployment, compute allocation, and biological integration of intelligent agents [21,23]. Third, ecological embedding of ethics should be prioritized: instead of aligning AI solely to human preferences, safety objectives must incorporate biospheric constraints, treating the stability of natural systems and interspecies coexistence as first-class alignment targets [24]. Finally, interpretability across substrates must become a central research agenda. As intelligence migrates into hybrid biological–computational forms, transparency cannot be limited to code or weights; it must extend to neural interfaces, biological feedback loops, and collective behavioral patterns [25].

Ultimately, the future challenge is not simply to prevent AI from escaping human control, but to decide what kind of civilization humans wish to inhabit once control itself becomes distributed. The coming era may require a shift from mastery to stewardship—from designing dominant systems to cultivating resilient cognitive ecosystems. Whether this transition leads to planetary instability or a new form of cooperative intelligence will depend less on any single technology than on humanity’s capacity to redefine its relationship with intelligence, life, and the Earth as a whole [3,11,12].



Figure 3: Four Future Directions. Autonomous agent networks and their integration with biological substrates challenge existing AI governance by undermining assumptions of centralized control and clear system boundaries. As intelligence becomes distributed across multi-agent and biohybrid ecosystems, the primary risk shifts from loss of oversight to loss of ontological clarity regarding agency and responsibility. This evolution forces a reconceptualization of humanity’s role within the biosphere, treating intelligence as a shared planetary resource rather than a human-controlled tool. Addressing these risks requires new strategies, including reversibility by design, global coordination mechanisms, ecological embedding of ethics, and cross-substrate

interpretability. Ultimately, AI safety must transition from technical control toward stewardship of complex cognitive ecosystems

Conclusion

The rapid emergence of autonomous agent networks and their increasing entanglement with biological substrates mark a fundamental transition in the trajectory of artificial intelligence. What began as isolated computational systems designed to optimize specific tasks is evolving into a distributed cognitive ecosystem, in which intelligence is no longer confined to discrete machines but embedded across technical infrastructures, social environments, and potentially living organisms. In this new landscape, traditional safety assumptions—centralized control, clear system boundaries, and reliable shutdown mechanisms—are progressively losing their practical and conceptual validity [26].

This shift compels a rethinking of AI not merely as a technological artifact, but as a systemic force that reshapes the conditions of agency, responsibility, and governance. The central challenge is no longer whether individual models can be aligned with human values, but whether complex, self-organizing intelligent networks can be governed under conditions of strategic interaction, ecological embedding, and partial biological integration. As intelligence becomes distributed and persistent, control becomes relational rather than hierarchical, and safety becomes an emergent property of institutions, incentives, and collective norms rather than a feature of code alone.

Ultimately, the future of AI safety will depend less on technical safeguards than on humanity’s capacity to construct new forms of coordination and stewardship. This includes designing systems with built-in reversibility, establishing global governance mechanisms, embedding ethical constraints at the level of ecosystems, and extending interpretability across heterogeneous substrates. Together, these shifts signal a transition from a paradigm of mastery to one of coexistence—from attempting to dominate intelligent systems to learning how to inhabit a shared cognitive environment. Whether this transition leads to instability or to a resilient form of planetary intelligence will depend on how successfully humans redefine their role within an emerging multi-agent, multi-species world.

Disclosure

Conflicts of Interest: Nothing to Disclose.

Ethical Approval

Not applicable.

Patient Consent

Not applicable.

Data Availability

Not applicable.

References

1. (2026) What is Moltbook? The strange new social media site for AI bots. The Guardian <https://www.theguardian.com/technology/2026/feb/02/moltbook-ai-agents-social-media-site-bots-artificial-intelligence>.
2. Varanasi L (2026) A look inside the Reddit-style social media site for AI agents that is dividing humans. Business Insider <https://www.businessinsider.com/moltbook-ai-agents-social-network-reddit-2026-2>.
3. Edwards B (2026) AI agents now have their own Reddit-style social network, and it’s getting weird fast. Ars Technica

- <https://arstechnica.com/information-technology/2026/01/ai-agents-now-have-their-own-reddit-style-social-network-and-its-getting-weird-fast/>.
4. Guo T, Chen X, Wang Y, Chang R, Pei S, et al. (2024) Large Language Model Based Multi-agents: A Survey of Progress and Challenges. IJCAI (Proceedings) <https://www.ijcai.org/proceedings/2024/890>.
 5. Wang L, Ma C, Feng X, Zhang Z, Yang H, et al. (2023) A Survey on Large Language Model based Autonomous Agents. arXiv <https://arxiv.org/pdf/2308.11432>.
 6. Brundage M, Avin S, Clark J, Toner H, Eckersley P, et al. (2018) The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation. arXiv <https://arxiv.org/pdf/1802.07228>.
 7. Hammond L, Oesterheld C, Reuel A, Chan A, Conitzer V, et al. (2025) Multi-Agent Risks from Advanced AI. arXiv <https://arxiv.org/pdf/2502.14143>.
 8. Ashery AF, Aiello LM, Baronchelli A (2025) Emergent social conventions and collective bias in LLM populations. *Science Advances* 11: eadu9368.
 9. Hadfield-Menell D, Dragan A, Abbeel P, Russell S (2016) Cooperative Inverse Reinforcement Learning. arXiv <https://arxiv.org/pdf/1606.03137>.
 10. Amodei D, Olah C, Steinhardt J, Christiano P, Schulman J, Mané D (2016) Concrete Problems in AI Safety. arXiv <https://arxiv.org/pdf/1606.06565>.
 11. (2023) AI Extinction Statement – Press Release. Center for AI Safety <https://safe.ai/work/press-release-ai-risk>.
 12. Perrigo B (2023) AI Is as Risky as Pandemics and Nuclear War, Top CEOs Say, Urging Global Cooperation. *TIME* <https://time.com/6283386/ai-risk-openai-deepmind-letter/>.
 13. Huang X, Liu W, Chen X, Wang X, Wang H, Lian D, et al. (2024) Understanding the planning of LLM agents: A survey. arXiv <https://arxiv.org/pdf/2402.02716>.
 14. Yang Y, Chai H, Song Y, Qi S, Wen M, et al. (2025) A Survey of AI Agent Protocols. arXiv <https://arxiv.org/pdf/2504.16736>.
 15. Carey R, Everitt T (2023) Human Control: Definitions and Algorithms (shutdown instructability / corrigibility variant). arXiv <https://arxiv.org/pdf/2305.19861>.
 16. Hudson R (2025) Corrigibility Transformation: Constructing Goals That Accept Updates[J]. arXiv <https://arxiv.org/pdf/2510.15395>.
 17. (2023) Artificial Intelligence Risk Management Framework (AI RMF 1.0). NIST <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>.
 18. (2024) Regulation (EU) 2024/1689 (AI Act). European Union <https://eur-lex.europa.eu/eli/reg/2024/1689/oj>.
 19. Thornley E (2024) The Shutdown Problem: An AI Engineering Puzzle for Decision Theorists. *Philosophical Studies* 182: 1653-1680.
 20. (2024) International AI treaty (“AI Convention”) adopted May 2024 (first legally binding international AI treaty). Council of Europe / Reuters.
 21. Smirnova L, Caffo SB, Gracias DH, Huang Q, Morales Pantoja IE, et al. (2023) Organoid intelligence (OI): the new frontier in biocomputing and intelligence-in-a-dish. *Frontiers in Science* 1:1017235.
 22. Hartung T, Morales Pantoja IE, Smirnova L (2024) Brain organoids and organoid intelligence from ethical, legal and social perspectives. *Front. Artif. Intell* 6: 1307613.
 23. Hettick M, Ho E, Poole AJ, Monge M, Papageorgiou D, Takahashi K, et al. (2025) Minimally invasive implantation of scalable high-density cortical microelectrode arrays. *Nature Biomedical Engineering* <https://www.nature.com/articles/s41551-025-01501-w#citeas>.
 24. Lee AH, Lee J, Leung V, Larson L, Nurmikko A (2024) Patterned electrical brain stimulation by a wireless network of implanted silicon microchips. *Nature Communications* 15: 10093.
 25. Li J, Chen G, Li G, Xiao L, Jia R, et al. (2025) Flexible brain electronic sensors advance wearable brain–computer interface technology (review). *NPJ Biomed Innov* 2: 24.
 26. Russell SJ (2019) *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking https://maxkasy.github.io/home/files/other/ML_Econ_Oxford/human_compatible.pdf.

Copyright: ©2026 Leilei Cheng, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.