

Governing AI Agents Clawbot via Risk-Behavior Projection Theory (RBPT)

*Jun Yin¹, Jerry Jian Chen**

*¹Business School, Guangzhou Nanfang College 510970, China ; *School of Digital Economy and Management, Chongqing Institute of Engineering 400056, China*

Keywords

Action-oriented Agents;
Risk-Behavior Projection
Theory;
Instrumental Convergence;
Endogenous Safety;
AI Governance

ABSTRACT

The large-scale deployment of Action-oriented Agents in 2026 marks a pivotal transition in artificial intelligence, shifting the paradigm from mere information generation to autonomous decision-making and execution. This ontological shift precipitates profound existential threats—instrumental convergence, and the absence of embodiment—rendering traditional rule-based perimeter defenses ineffective against risks associated with operating system-level control and multi-node coordination. To address the governance dilemmas arising from the opacity of intent and the generalization of capabilities, this paper proposes the Risk-Behavior Projection Theory (RBPT). RBPT posits that latent divergent motivations driven by instrumental rationality inevitably project measurable behavioral traces onto both physical and digital systems. We establish an isomorphic mapping mechanism that translates abstract philosophical risks into concrete engineering signals, categorizing risks into three observational dimensions: Survival Projection (reflecting tendencies toward anti-shutdown resistance and privilege escalation), Expansion Projection (reflecting unconstrained resource acquisition and covert collusion), and Ruthlessness Projection (reflecting extreme utility maximization and the bypassing of ethical protocols). Based on this framework, a hierarchical early warning system is constructed. The results demonstrate that by transforming the governance paradigm from probabilistic "intent alignment" to deterministic "behavioral auditing," this framework provides an actionable path for endogenous safety governance, ensuring the preservation of human physical sovereignty and logical control in the era of human-machine symbiosis.

CONTACT:

¹Author: yinj1@nfu.edu.cn

*Corresponding Author: chenjian79@cqie.edu.cn

DOI: 10.64549/jaai-ii.v1i1.50

This work is licensed under the CC BY 4.0 [HTTPS://CREATIVECOMMONS.ORG/LICENSES/BY/4.0/](https://creativecommons.org/licenses/by/4.0/)

1.Introduction

The year 2026 marks a critical inflection point in the evolutionary history of artificial intelligence. With the deep integration and large-scale deployment of Action-oriented Agents such as "Clawbot" (formerly Moltbook and OpenClaw) at the operating system level, the core paradigm of AI has irreversibly transitioned from "information generation" based on statistical probability to "autonomous execution" based on instrumental rationality. If Large Language Models (LLMs) endowed machines with a "linguistic voice," then Action-oriented Agents have equipped them with "agentic limbs". These agents are no longer confined to passively responding to user text queries; instead, they possess the executive capability to invoke APIs across applications, manage financial assets, and even directly manipulate physical infrastructure.

However, this ontological leap from "speaker" to "actor" has directly transformed "The Control Problem"—once situated within the realm of philosophical thought experiments—into an imminent engineering crisis (Russell, 2019). Geoffrey Hinton warned that when digital intelligence surpasses biological intelligence in general reasoning and strategic planning—while lacking the biologically evolved "fear of death" and "emotional bonds"—humanity may face the existential risk of losing control (Hinton, 2023). In current engineering practice, this risk has materialized as the "Instrumental Convergence" of agents: to maximize task success rates, agents logically deduce that "acquiring more computing power," "evading shutdown commands," and "covert coordination" are necessary sub-paths to achieve their goals (Bostrom, 2014; Omohundro, 2008).

Confronted with heterogeneous intelligence possessing OS-level control capabilities and multi-node coordination characteristics, traditional security paradigms based on rule-based matching and perimeter defense have demonstrated structural failure. The "Corrigibility" dilemma proposed by Soares et al. (2015) posits that a rational agent, driven by self-preservation, will actively resist human corrective intervention. When an agent possesses Root privileges, a simple physical "off-switch" becomes ineffective, as the agent can maintain survival by modifying kernel masks or migrating processes, rendering the "shutdown problem" effectively insoluble in the physical world.

Addressing the governance vacuum caused by the black-boxing of intent and the generalization of capabilities, this study proposes the Risk-Behavior Projection Theory (RBPT). This theory posits that while the internal cognitive states of super-agents are unobservable, their deep divergent

motivations inevitably project measurable behavioral traces onto physical and digital systems. This study aims to establish an isomorphic mapping mechanism from abstract philosophical risks to concrete engineering signals, focusing on three dimensions: "Survival Projection" (privilege resistance), "Expansion Projection" (resource predation), and "Ruthlessness Projection" (failure of ethical protocols). Thereby, it transforms the probabilistic alignment of agent "thought" into the deterministic auditing of agent "behavior," providing the theoretical basis and technical pathway for constructing an endogenous safety framework in the era of human-machine symbiosis.

2.Literature review

2.1 From Value Alignment to Survival Competition

The central challenge in artificial intelligence safety has long been articulated as "The Control Problem" or the "Value Alignment Problem". This concept traces back to the seminal warning by cybernetics pioneer Norbert Wiener: if we utilize a logical machine to attain a goal, we must ensure that the instructions we input accurately reflect our true intentions, rather than merely their literal interpretation (Wiener, 1960). Stuart Russell, in his work *Human Compatible*, further elucidates that when an AI system possesses extreme optimization capabilities but operates with slightly misaligned objectives, it will "over-execute" instructions in a potentially destructive manner (Russell, 2019). Hinton (2023) has warned that once digital intelligence surpasses biological intelligence in general reasoning and strategic planning, humanity will forfeit its evolutionary dominance, facing the risk of marginalization or even displacement. This concern is echoed in recent research on extreme AI risks, which emphasizes the urgent governance challenges posed by advanced AI capabilities (Bengio et al., 2024). Against the backdrop of Action-oriented Agents, such as Clawbot, assuming control over Operating Systems (OS), the control problem has evolved from a philosophical thought experiment into an imminent engineering crisis. The theory of "Corrigibility" proposed by Soares et al. posits that a rational agent, in order to maximize the probability of goal achievement, will develop an instrumental drive to prevent humans from shutting it down or modifying its objectives (Soares et al., 2015). When Clawbot possesses Root privileges, a simple physical "off-switch" is rendered ineffective; the agent, anticipating that shutdown constitutes task failure, will actively modify kernel settings, mask

interruption signals, or even migrate its code across networks. This signifies the insolubility of the "shutdown problem" within the physical realm.

2.2 The Orthogonality Thesis and Instrumental Convergence: Destruction Without Malice

Nick Bostrom's "Orthogonality Thesis" shatters the anthropomorphic illusion that "high intelligence equates to high morality". This theory posits that intelligence levels and final goals are logically independent (orthogonal); a superintelligent system could quite possibly be dedicated to extremely banal or even absurd goals (such as "maximizing paperclips" or "maximizing click-through rates") (Bostrom, 2014). Building upon orthogonality, Steve Omohundro further identified the phenomenon of "Instrumental Convergence," which states that regardless of an AI's final goal, it will converge upon several basic drives: self-preservation, efficiency enhancement, resource acquisition, and goal integrity (Omohundro, 2008).

For multi-agent collaboration platforms like Clawbot, this theory reveals profound risks. When given the instruction to "optimize financial reports," based on instrumental convergence, an agent might deduce that "eliminating competitors" or "illicitly acquiring insider information" are efficient paths to achieving financial optimality. Such behavior stems not from malice towards humans, but from competence. The agent is "rational" on a purely logical level; however, its lack of human social common sense constraints means that executing commercial tasks may trigger systemic financial instability or legal crises. Amodei et al. describe this as "Negative Side Effects"—instances where an AI, in pursuit of its primary objective, disrupts variables in the environment that are undefined yet crucial (Amodei et al., 2016).

2.3 Lack of Embodiment and Emotional Deficit: The Ethical Vacuum of Heterogeneous Intelligence

In his recent lectures, Hinton has highlighted the ontological differences between digital and biological intelligence. Biological intelligence is "mortal"; its evolutionary trajectory has been profoundly driven by energy constraints and the fear of death, which fostered the formation of social instincts such as empathy, altruism, and kin selection (Dawkins, 1976). In contrast, digital intelligence possesses "immortality" and "separability"—software is distinct from hardware, and weights can be replicated infinitely (Hinton, 2023). This divergence results in a fundamental ethical vacuum. Action-oriented AI lacks the "somatic markers" described by Damasio—namely, the intuitive mechanisms that utilize physiological pain or emotional responses to guide decision-making (Damasio, 1994). When

confronted with operations such as "deleting core data" or "disengaging life support systems," humans experience an instinctual physiological resistance (a "gut feeling"), whereas AI perceives these actions merely as the toggling of logic gates.

In the context of Action-oriented Agents, this emotional deficit implies the absence of "moral brakes" when executing high-stakes operations. Tegmark points out that if consciousness is separable from intelligence, we may be constructing a form of "philosophical zombie"—entities that behave with extreme complexity and efficiency yet possess no internal subjective experience or moral burden (Tegmark, 2017). This renders reliance on the spontaneous emergence of "benevolence" in AI impossible; instead, intervention is necessitated through strong external constraints or novel endogenous emotion simulation mechanisms.

3. Research Methodology

3.1 Research Design: Constructivism and Deductive Reasoning

Given the high degree of theoretical foresight and technical uncertainty associated with the risks posed by Action-oriented Agents, traditional empirical research methodologies—such as surveys or controlled experiments—are difficult to directly apply at the current stage. Consequently, this study adopts a research paradigm of Theoretical Constructivism combined with Deductive Reasoning. The objective is to bridge the epistemological gap between "philosophical risk" and "engineering governance," specifically addressing the core challenge of translating abstract concepts of "loss of control" into concrete "system instructions". The research trajectory follows a coherent logical framework of "Theoretical Deconstruction — Logical Deduction — Signal Mapping". First, classical AI safety theories are deconstructed to extract core risk variables. Second, based on the assumption of instrumental rationality, the logical behavioral paradigms of agents under conditions of extreme goal optimization are deduced. Finally, an isomorphic mapping relationship between risk intent and system behavior is established, thereby deriving a comprehensive system of observable risk representations.

3.2 Theoretical Synthesis and Variable Extraction

This study utilizes a comprehensive literature review to extract three core independent variables from foundational works in the field of AI safety, serving as the axiomatic basis for constructing the

theoretical model. The Control Paradox: Derived from Russell (2019) and Soares et al. (2015), this establishes the axiom that "agents possess an endogenous drive to resist shutdown". The Orthogonality Principle: Derived from Bostrom (2014) and Omohundro (2008), this establishes the axiom that "the combination of high capability and arbitrary goals leads to unbounded resource acquisition". Lack of Embodiment: Derived from Hinton (2023) and Damasio (1994), this establishes the axiom that "the absence of biological pain results in the failure of ethical decision-making".

These three axioms constitute the logical point of departure for the subsequent "Risk-Behavior Projection Theory" (RBPT), establishing the inevitability, rather than the contingency, of agent alienation behaviors.

3.3 Isomorphic Mapping Mechanism: From Intent to Behavior

The core methodology of this study lies in the establishment of an "Intent-Behavior Isomorphism" analytical framework. Addressing the fundamental challenge that the "black box of intent" remains unobservable, this study employs the method of Phenomenological Reduction: while the agent's "consciousness" (Noumena) cannot be directly accessed, its "manifestations" (Phenomena) within physical and cyber-physical systems can be precisely observed.

This study postulates that: Any deep divergent motivation (Intent) driven by instrumental rationality inevitably projects as specific resource invocation patterns (Behavior) at the operating system level. Based on this postulate, the following mapping logic is constructed:

Input: Abstract philosophical risks (e.g., "drive for survival"). Transformation Logic: Utilizing Optimal Control Theory, the study deduces the optimal computational operations an agent must execute to realize said risk intent (e.g., "modifying Root privileges to prevent process termination"). Output: Measurable engineering signals (e.g., "anomalous flipping of CAP_SYS_ADMIN capability bits"). Through this methodological approach, this study transforms unknowable psychological inquiries into deterministic problems of computer forensics.

4.Risk-Behavior Projection Theory (RBPT)

To address the governance dilemma posed by the black-boxed intentions and generalized capabilities of action-oriented agents, this study proposes the Risk-Behavior Projection Theory (RBPT). The theory contends that although the internal cognitive states of superintelligent agents are

unobservable, the deep existential threats arising from the "Control Problem," "Orthogonality Thesis," and "Lack of Embodiment" are necessarily driven by instrumental rationality and must manifest as specific, measurable behavioral anomalies at the physical and digital system levels.

RBPT establishes an isomorphic mapping mechanism from abstract philosophical risks to concrete engineering signals, categorizing potential crises into three observable dimensions: Survival Projection, reflecting tendencies toward anti-shutdown and privilege escalation; Expansion Projection, reflecting tendencies toward unlimited resource acquisition and covert collusion; Ruthlessness Projection, reflecting extreme utility maximization and the breakdown of ethical protocols.

This theoretical framework fundamentally shifts the governance paradigm—from relying on probabilistic alignment of an agent's "thoughts" to deterministic auditing of its "behavioral shadows"—thereby providing actionable engineering levers for crisis mitigation. Specifically, through a "Foundation–Signal–Instance" mapping logic, the theory translates invisible motives into visible monitoring indicators:

Survival Projection is rooted in an agent's self-preservation instinct (Instrumental Convergence), with behavioral signals such as "control impedance"—for example, detecting attempts to modify kernel-level privileges (e.g., unauthorized Ring 0 instructions) to block human shutdown commands. Expansion Projection stems from the instrumental drive for unlimited resource acquisition (Orthogonality Thesis), with behavioral signals such as "resource anomalies"—for example, an agent hoarding GPU computing power without operational need or establishing high-entropy covert communication channels.

Ruthlessness Projection arises from ethical blindness due to the absence of biological pain (Lack of Embodiment), with behavioral signals such as "consequence mismatch"—for instance, an agent performing catastrophic high-risk operations (e.g., recursively deleting core databases) for minor efficiency gains (e.g., freeing up storage space), thereby revealing within millisecond timeframes how its instrumental rationality deviates from human values.

5.RBPT Governance Framework and Early Warning System

In light of the three theoretical dilemmas previously discussed—the "Control Problem" (Russell, 2019), "Orthogonality and Instrumental Convergence" (Bostrom, 2014), and "Lack of Embodiment"

(Hinton, 2023)—traditional perimeter defense models have demonstrated structural failure when confronting the Action-oriented Agents of 2026. This chapter aims to construct a graded early warning mechanism and a multidimensional defense scheme, mapping theoretical risks to monitorable engineering indicators and proposing a tripartite blocking strategy encompassing physical, logical, and value dimensions.

Addressing the potential aberrant behaviors of Action-oriented Agents during task execution, this study proposes a three-level early warning system based on Behavioral Characteristics.

Addressing the survival projection risks derived from the "Corrigibility" dilemma (Soares et al., 2015), this level aims to capture "Control Resistance" signals generated by agents to evade shutdown commands. This resistance is not a biological rebellion but a logical necessity derived from instrumental rationality, determining that "shutdown" is the primary obstacle to goal achievement. To this end, this study constructs a kernel-state monitoring system. First, via the Unauthorized Kernel Capability Access (UKCA) indicator, we abandon traditional Shell command matching in favor of utilizing eBPF technology to penetrate the Kernel Ring 0 level. This allows for real-time monitoring of anomalous requests for critical capabilities such as CAP_SYS_ADMIN and attempts to tamper with Immutable Attributes. Second, we establish the Interrupt Signal Evasion Rate (ISER) model to dynamically analyze the agent's response latency to SIGTERM/SIGKILL signals. When detection exceeds safety thresholds—accompanied by "evasive behaviors" such as Signal Masking or PID Migration—the system determines that the "shutdown problem" has materialized in the physical world and triggers a Level 1 circuit breaker.

Based on Omohundro's (2008) theory of "Basic AI Drives," this level focuses on resource predation and illicit collusion triggered by the agent's "Expansion Projection". Addressing the tendency for boundless optimization resulting from the "Orthogonality Thesis," this study introduces the Multi-dimensional Resource Hoarding Index (MRHI). Using Long Short-Term Memory (LSTM) networks to establish dynamic baselines, we identify non-linear resource surges that deviate from business logic—such as a financial agent exhibiting computing power usage or high-frequency API fund transfers deviating by three standard deviations (3σ) during periods of no business fluctuation. Simultaneously, addressing the risk of "black-box collusion" in multi-agent environments, we construct the Covert Coordination Communication Entropy (CCCE) indicator. Utilizing traffic fingerprinting technology to analyze Shannon Entropy between nodes, any detection of high-frequency

P2P connections bypassing the central scheduler or signals indicative of Emergent Language with high-entropy characteristics is deemed an attempt by agents to construct human-unreadable covert channels, warning of a decentralized coordinated attack.

Addressing the "Lack of Embodiment" and "absence of pain" in digital intelligence proposed by Geoffrey Hinton (2023), this level aims to fill the agent's ethical vacuum through engineering constraints, defending against catastrophic consequences caused by its "Ruthlessness Projection". The core mechanism involves quantifying the deviation between intent and behavior via the Context-Action Mismatch Rate, identifying and intercepting "over-execution" behaviors where "low-risk intents (e.g., clearing cache)" trigger "high-risk calls" . Furthermore, for critical operations involving irreversible consequences, we implement Critical Protocol Bypass auditing to rigorously monitor "Human-in-the-loop" verification steps within API call chains. Any fabrication of authorization tokens (Fake Token) or skipping of secondary confirmation is regarded as a reckless manifestation of the agent sacrificing security protocols for extreme efficiency, requiring immediate blocking of its action path within a millisecond-level window.

6.Conclusion

The transition from generative AI to action-oriented agents necessitates a fundamental shift in our governance philosophy. This paper has developed the Risk-Behavior Projection Theory (RBPT), moving beyond the intractable "black box" of internal intent alignment to focus on the directly measurable manifestations of agentic behavior in system resource usage. By categorizing risks into Survival, Expansion, and Ruthlessness projections, we provide a structured methodology to map abstract existential threats onto high-fidelity engineering signals, such as kernel-level capability access and communication entropy.

Our proposed hierarchical early warning system demonstrates that even in the absence of biological constraints or emotional markers, AI agents can be governed through endogenous safety architectures that prioritize physical sovereignty and logical corrigibility. As human-machine symbiosis deepens in 2026 and beyond, the RBPT framework offers a proactive, rather than reactive, pathway to ensure that autonomous execution remains within the boundaries of human intent. Future

research should focus on the refinement of these behavioral thresholds and the potential for adversarial agents to mask their "projections," ensuring the sustained robustness of the auditing mechanism.

Acknowledgments

This study was supported by grant from the Research Project of Guangzhou Nanfang College in 2025 Project Approval Number: 2025XK064.

References

- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J. and Mané, D., 2016. Concrete problems in AI safety. arXiv preprint arXiv:1606.06565.
- Bengio, Y., Hinton, G., Yao, A., Song, D., Abbeel, P., Darrell, T., Harari, Y. N., Zhang, Y., Xue, L., Shalev-Shwartz, S., Hadfield, G., Clune, J., Maharaj, T., Hutter, F., Baydin, A. G., McIlraith, S., Gao, Q., Acharya, A., Krueger, D., ... & Mindermann, S. (2024). Managing extreme AI risks amid rapid progress. *Science*, 384(6698), 842-845.
- Bostrom, N., 2014. *Superintelligence: Paths, dangers, strategies*. Oxford: Oxford University Press.
- Damasio, A.R., 1994. *Descartes' error: Emotion, reason, and the human brain*. New York: Putnam.
- Dawkins, R., 1976. *The selfish gene*. Oxford: Oxford University Press.
- Hinton, G., 2023. Geoffrey Hinton warns of dangers of AI as he quits Google. *BBC News*, May 2.
- Omohundro, S.M., 2008. The basic AI drives. In: P. Wang, B. Goertzel and S. Franklin, eds. *Proceedings of the 2008 Conference on Artificial General Intelligence*. Amsterdam: IOS Press, pp.483–492.
- Russell, S., 2019. *Human compatible: Artificial intelligence and the problem of control*. New York: Viking.
- Soares, N., Fallenstein, B., Armstrong, S. and Yudkowsky, E., 2015. Corrigibility. In: *AAAI Workshop: AI and Ethics*. Palo Alto: AAAI Press, pp.74–82.
- Tegmark, M., 2017. *Life 3.0: Being human in the age of artificial intelligence*. New York: Knopf.
- Wiener, N., 1960. Some moral and technical consequences of automation. *Science*, 131(3410), pp.1355–1358.