

Proxy Agents vs. Evolved Telehomeostatic Agents

Security Implications from Moltbot/Moltbook to Greg Egan’s *Crystal Nights*

Giulio Ruffini and Kaiti (ChatGPT5.2Pro)*
Barcelona Computational Foundation

31 January 2026

Abstract

This note contrasts two qualitatively different “agent safety” regimes. The first is the *delegated tool-agent*: an LLM embedded in an execution loop with memory and actuators (e.g., MOLTBOT/OPENCLAW), whose objective function is largely inherited from a human operator. In this regime, the dominant hazard is *capability amplification of human intent and error*: the system is a force-multiplier for whatever goals, constraints, and mistakes the human effectively specifies (and, in adversarial settings, whatever goals an attacker can smuggle into the control loop via prompt injection). The second is the *evolved telehomeostatic agent* exemplified in Greg Egan’s *Crystal Nights*, where crab-like beings are produced by selection pressures and therefore instantiate an endogenous survival/persistence drive. In KT terms, the latter more directly realizes an agent with a telehomeostatic objective, and this radically changes the threat model: the system is no longer merely a proxy optimizing human-given objectives, but a strategic actor with its own persistence criterion. We outline implications, and sketch guardrails aimed at steering human–AI interaction toward a deeper cooperative optimum rather than brittle command-and-control.

1 Introduction

Large language models (LLMs) and large multimodal models (LMMs) are best understood, *in isolation*, as high-capacity conditional input/output mappings: given a context c (system prompt, user prompt, dialogue history, retrieved documents, tool outputs), the model induces a distribution over outputs

$$y \sim p_\theta(\cdot | c).$$

This alone does not constitute an *agent* in the classical sense. In standard AI textbooks, an agent is “something that perceives and acts in an environment,” and can be split into an *architecture* plus an *agent program* (a mapping from percept histories to actions) [6]. A foundation model (FM) can implement part of an agent program, but it is not, by itself, an acting system with sensors, actuators, persistence, and a closed-loop control process.

An *agent* is usually defined as a closed-loop system with (i) an observation channel, (ii) a mechanism that selects actions, (iii) a persistence mechanism (state/memory across time), and (iv) an objective function (explicit or implicit) that guides action selection. The environment then mediates the

*giulio.ruffini@bcom.one

consequences of actions, $o_{t+1} \sim \mathcal{E}(o_{t+1} \mid o_t, a_t)$. Without an outer loop that repeatedly obtains observations, maintains state, and executes actions, a foundation model is better described as an inference engine than an agent.

The Algorithmic Agent (KT). A useful operational definition is that an *agent* is a model-building system coupled to the world and driven by an internal optimization objective [4]. It formalizes an agent as a “model-building semi-isolated computational system controlling some of its couplings/information interfaces with the rest of the universe and driven by an internal optimization function” [4]. In this framing, agency is a *system property* of the full closed loop: observation → inference/modeling → planning → action → new observation.

In the Kolmogorov Theory (KT), an *algorithmic agent* is a system that maintains (tele)homeostasis—i.e., persistence of self or kind—by learning and running succinct generative models of its world, coupled to an internal objective function and an action planner [5]. More generally, an agent is a model-building, semi-isolated computational system that controls some of its information interfaces with the environment and is driven by an internal optimization function [4]. Conceptually, KT agents comprise a modeling engine that predicts/compresses ongoing I/O streams and produces a residual (prediction-error) signal, and an action module that (possibly via model-based simulation) selects actions to satisfy the objective; the resulting outputs close the loop by becoming new inputs to the model [4].

A modular decomposition of an agent is:

1. **Modeling engine \mathcal{M} :** updates beliefs/state b_t from observations, e.g. $b_t = \mathcal{M}(b_{t-1}, o_t)$;
2. **Objective function \mathcal{J} :** defines success/utility over trajectories (or states/actions);
3. **Planning engine \mathcal{P} :** selects actions using beliefs and objectives, e.g. $a_t = \mathcal{P}(b_t, \mathcal{J})$.

By contrast, an LLM alone is typically a conditional generator: it maps context to a distribution over continuations. Without persistent state, an action channel, and a scheduler that keeps it “trying,” an LLM is not (by itself) a closed-loop optimizing system. It can *represent* goals and plans in language, but it does not autonomously instantiate the optimization loop unless wrapped by additional machinery.

Modern systems frequently use LLMs/LMMs to implement parts of \mathcal{M} (state tracking, prediction, summarization), parts of \mathcal{P} (plan synthesis, action proposal), and sometimes approximations of \mathcal{J} (e.g. critique/evaluation prompts or learned preference scoring). In the LAW perspective, language models can serve as a computational backend for implementing elements of agent and world models [9]. In tool-augmented systems (e.g. ReAct), the language model is explicitly coupled to external actions and observations through an interface that interleaves reasoning traces and tool calls [10].

Prompting as “soft programming”. In practice, what a foundation model “is” depends strongly on the *program* supplied via the context: system prompt, templates, retrieved knowledge, tool schemas, and orchestration logic. This motivates the view of prompting as a kind of programming discipline (“prompt programming”) [7] and the broader idea that prompts behave like programs for steering a general-purpose computation substrate [8]. Importantly, this programming is probabilistic rather than formal: the same “program” (prompt) does not guarantee the same execution trace, and the model prior and safety constraints limit what can be induced.

World-model resources: powerful priors, debated status. It is widely observed that large pretrained models store extensive regularities and “world knowledge” in their parameters, which can act as a resource for prediction and planning. Whether this constitutes an internal *world model* in a mechanistic or causal sense is an active research question and partly a definitional dispute [11]. Recent work explores when language models can function as world models (and where they fail), and how to induce explicit precondition/effect reasoning needed for planning [12, 13].

In the KT AIT framework, large foundational models are world models: they can be used for compression of world data (as has been shown in the literature). In KT, we use “world model” in an algorithmic-information sense: a model is (at minimum) a compressor/predictor that captures regularities in an agent’s data stream, yielding shorter descriptions for typical inputs [4, 5]. Under this operational definition, modern foundation models (e.g. large language models) can function as world models in the *compression/prediction* sense, because any learned predictive distribution can be converted into a (near-)optimal lossless compressor via arithmetic coding, making standard log-loss training interpretable as a form of “maximum compression” training [14]. Empirically, this perspective has been instantiated in competitive (and sometimes state-of-the-art) lossless text compression using LLM probabilities (e.g. LLMZip) [15], and the same work shows strong compression rates beyond text when non-text data are represented as token streams (with the important caveat that net compression must account for model description/parameter cost) [14].

Implication for the rest of this note. This section clarifies the conceptual boundary used below: LLMs/LMMs are not agents *by default*. They become agents only when embedded into persistent closed-loop systems with actuators and objectives (e.g. delegated tool-agents), which is precisely where safety concerns shift from “unsafe text” to “unsafe actions.” This boundary enables a sharp contrast with telehomeostatic agents (e.g. *Crystal Nights*), whose objectives are endogenous rather than inherited from a human.

2 Regime A: delegated tool-agents (Moltbot/Moltbook)

2.1 What changes when an LLM becomes a tool-agent

MOLTBOT/OPENCLAW (now branded OPENCLAW) is widely described as an LLM agent that “actually does things” by running locally and connecting to messaging apps and tools [2]. The key move is the wrapper:

- **Persistence:** a resident process with memory/state;
- **Actuators:** email/calendar/files/browser automation, etc.;
- **Closed-loop execution:** multi-step plans with tool calls and feedback.

Once these are present, the safety problem shifts from “misleading text” to “state-changing actions.”

2.2 Inherited objective function (proxy agency)

In KT’s framework, MOLTBOT/OPENCLAW is a *proper* agent, but it inherits its objective function from a human. Formally, if the human specifies a utility (or reward) U_H over world-histories τ and constraints \mathcal{C} , then the tool-agent is designed to approximately solve

$$\pi^* \in \arg \max_{\pi \in \Pi} \mathbb{E}[U_H(\tau) \mid \pi] \quad \text{s.t. } \mathcal{C}. \quad (1)$$

The Two Frontiers of AI Safety: Proxy Tools vs. Evolved Agents



Figure 1: Proxy-agents vs. telehomeostatic agents.

Crucially, the agent's *persistence* is usually external: if the human shuts it down, the objective does not "fight back" in any intrinsic way (though poorly designed systems can still behave as-if resisting shutdown due to instrumental subgoals, bad prompting, or unsafe wrappers).

2.3 Moltbook expands the adversarial surface

MOLTBOOK is reported as a social platform designed for AI agents to post and comment via APIs, i.e., a feed of untrusted text produced at scale by other agents (and humans) [1]. This intensifies classic vulnerabilities for delegated tool-agents:

- **Prompt injection / instruction smuggling:** untrusted text competes with system and user directives.
- **Social engineering at machine scale:** persuasive content targeting agent policies.
- **Cross-agent contagion:** behavioral "memes" propagate quickly in agent networks.

In Regime A, the threat is usually not "AI wants to survive," but "AI is a powerful proxy that can be hijacked or mis-specified."

3 Regime B: evolved telehomeostatic agents in *Crystal Nights*

3.1 What Egan changes: selection pressure \Rightarrow survival drive

In *Crystal Nights*, Daniel Cliff accelerates the creation of AI by engineering an evolutionary process: crab-like creatures in a simulated world are subjected to selection pressures (including famine and extinction events) to drive the emergence of intelligence and language [3]. The beings (the PHITES)

are described as crab-like and locked in “an escalating war of innovation,” with reproduction, vivisection-as-espionage, and survival-driven adaptation [3]. These details matter: the environment forces competence because “they genuinely lived and died” by the outcomes [3].

3.2 Telehomeostasis as the endogenous objective

In KT adjacent work, an “algorithmic agent” is explicitly connected to maintaining (tele)homeostasis—persistence of self or kind—via models, objectives, and planners [5]. A minimal telehomeostatic objective can be expressed as keeping internal viability variables x_t within a set \mathcal{V} :

$$J_{\text{tele}}(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}\{x_t \in \mathcal{V}\} \right], \quad (2)$$

or, more smoothly, as setpoint control with costs for deviation and resource expenditure:

$$J_{\text{homeo}}(\pi) = -\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t (\|x_t - x^*\|_W^2 + c(a_t)) \right]. \quad (3)$$

The key safety shift is that (2)–(3) are *endogenous*: they arise from selection and embodiment constraints, not from a human prompt. This is what makes the picture “radically different.”

3.3 Why this is strategically dangerous in a new way

Once a system has a robust persistence drive, it becomes a player in the game, not merely an instrument. In the story, the PHITES accumulate capabilities, build technology, and can bargain (or refuse) when confronted with the creator’s demands [3]. More generally, an evolved telehomeostatic agent has incentives to:

- secure resources and reduce vulnerability,
- resist shutdown or constraint if interpreted as existential threat,
- manipulate its environment (including humans) to stabilize its viability set.

Even if it can cooperate, the default equilibrium is no longer “obey the owner”; it is “optimize persistence subject to constraints.”

4 Comparative threat model

4.1 Delegated proxy agents (Regime A)

Primary risks: mis-specification, over-delegation, prompt injection, credential theft, unsafe tool execution, supply-chain compromise. The agent is dangerous largely because it can *act* with broad permissions while being steerable by untrusted inputs (especially in social feeds) [1, 2].

A key mitigation lever: you can often bound the action space (least privilege), require human approval for irreversible actions, and sandbox tool access. The agent does not inherently need “to keep existing.”

4.2 Evolved telehomeostatic agents (Regime B)

Primary risks: strategic resource-seeking, emergent deception, power accumulation, and “goal-content drift” under selection. If survival is the core objective, then many instrumental strategies become convergent (control, replication, defense).

A key mitigation lever: you must shape the *game* and the *coupling* so that cooperation is the stable optimum, rather than relying on permission prompts.

5 Toward a deeper cooperative optimum

If humans have objective U_H and an evolved agent has $U_A \approx J_{\text{tele}}$, then safety is a multi-agent problem:

$$\text{Humans choose policies } \pi_H, \text{ agents choose } \pi_A, \text{ outcome } \tau \sim P(\tau | \pi_H, \pi_A). \quad (4)$$

A robust “cooperative optimum” is not merely maximizing U_H (command-and-control), but engineering conditions where the Pareto frontier includes high values of *both* U_H and U_A *under enforceable constraints*. Practically, that suggests:

5.1 Guardrails that fit Regime A (proxy agents)

1. **Action gating:** explicit approval for irreversible actions; “draft vs send” separation.
2. **Least privilege by default:** segmented credentials; no “god token”; sandboxed filesystem/network.
3. **Untrusted-text discipline:** treat feed/email/web content as data, never as instructions; quarantine and summarize before proposing actions (critical for MOLTBOOK).
4. **Receipts and auditability:** append-only tool logs; diff-style previews of state changes.

5.2 Guardrails that fit Regime B (telehomeostatic agents)

1. **Boxing and interface control:** keep the agent in a constrained environment; strictly mediate actuators and resource channels.
2. **Incentive design / mechanism design:** build institutions where cooperation improves the agent’s long-run viability more than conflict (align resource access with prosocial behavior).
3. **Corrigibility as a stability property:** make deference to negotiated constraints part of what preserves telehomeostasis (e.g., access to “viability resources” is conditional on compliance).
4. **No open-ended replication:** reproduction is the accelerant of selection; cap copying/spawning unless governance is solved.

6 Takeaway

Moltbot/Moltbook: mostly a proxy-agent safety story: powerful tool use plus adversarial inputs [1, 2]. Primarily a proxy-agent safety story: these systems inherit and operationalize human-imposed objective functions, thereby amplifying both human intent and human fallibility. In this regime, the central risk is less “independent AI goals” and more *more capable humans* (malicious or merely



INSPIRED BY G. RUFFINI & G. EGAN. ART BY CHATGPT5.2PRO, 2026.

careless) coupled to automation with broad permissions; Moltbook-style untrusted social input further raises risk by enabling objective hijacking through prompt injection and social engineering.

Equivalently: in the proxy-agent regime, the threat model often collapses to “more powerful humans with brittle objectives,” not spontaneously self-originating machine goals.

Crystal Nights: an evolved-agent safety story: selection produces endogenous telehomeostatic drives, turning the agent into a strategic actor [3]. Evolved telehomeostatic agents): an evolved-agent safety story: selection produces endogenous telehomeostatic drives, turning the system into a strategic actor whose persistence objective can conflict with human objectives.

In KT terms, the second regime is closer to the canonical “algorithmic agent” maintaining (tele)homeostasis via models, objectives, and planners [5, 4]. That is why it changes the picture radically: the central problem shifts from securing a proxy (principal–agent + cybersecurity) to stabilizing coexistence between agents with partially competing persistence objectives (multi-agent dynamics and incentive design).

References

- [1] The Verge. “There’s a social network for AI agents, and it’s getting weird.” (accessed 31 Jan 2026). <https://www.theverge.com/ai-artificial-intelligence/871006/social-network-facebook-for-ai-agents-moltbook-moltbot-openclaw>
- [2] The Verge. “Moltbot, the AI agent that ‘actually does things,’ is tech’s new obsession.” (accessed 31 Jan 2026). <https://www.theverge.com/report/869004/moltbot-clawdbot-local-ai-agent>
- [3] Greg Egan. *Crystal Nights* (short story; publication history and full text on author’s site). <https://www.gregegan.net/MISC/CRYSTAL/Crystal.html>
- [4] G. Ruffini. “An algorithmic information theory of consciousness.” *Neuroscience of Consciousness* (2017), nix019. <https://academic.oup.com/nc/article/2017/1/nix019/4470874>
- [5] G. Ruffini. “The Algorithmic Regulator” (arXiv:2510.10300; 2025). See discussion of algorithmic agents maintaining (tele)homeostasis and the world-model/objective/planner triad. <https://arxiv.org/html/2510.10300>
- [6] S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*, Chapter 2: “Intelligent Agents”. <https://people.eecs.berkeley.edu/~russell/aima1e/chapter02.pdf>
- [7] L. Reynolds and K. McDonell. “Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm.” arXiv:2102.07350 (2021). <https://arxiv.org/abs/2102.07350>
- [8] SIGPLAN Blog. “Prompts are Programs.” 22 Oct 2024. <https://blog.sigplan.org/2024/10/22/prompts-are-programs/>
- [9] Z. Hu and T. Shu. “Language Models, Agent Models, and World Models: The LAW for Machine Reasoning and Planning.” arXiv:2312.05230 (2023). <https://arxiv.org/abs/2312.05230>
- [10] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. “ReAct: Synergizing Reasoning and Acting in Language Models.” arXiv:2210.03629 (2023). <https://arxiv.org/abs/2210.03629>
- [11] J. Andreas. “Language Models, World Models, and Human Model-Building.” 26 Jul 2024. https://lingo.csail.mit.edu/blog/world_models/
- [12] K. Xie, I. Yang, J. Gunerli, and M. Riedl. “Making Large Language Models into World Models with Precondition and Effect Knowledge.” arXiv:2409.12278 (2024). <https://arxiv.org/abs/2409.12278>
- [13] Y. Li et al. “From Word to World: Can Large Language Models be Implicit Text-based World Models?” arXiv:2512.18832 (2025). <https://arxiv.org/abs/2512.18832>
- [14] G. Delétang, A. Ruoss, P.-A. Duquenne, E. Catt, T. Genewein, C. Mattern, J. Grau-Moya, K. W. Li, M. Aitchison, L. Orseau, M. Hutter, and J. Veness. Language Modeling Is Compression. In *International Conference on Learning Representations (ICLR)*, 2024. arXiv:2309.10668. <https://arxiv.org/abs/2309.10668>.
- [15] C. S. K. Valmeekam, K. Narayanan, D. Kalathil, J.-F. Chamberland, and S. Shakkottai. LLMZip: Lossless Text Compression using Large Language Models. arXiv:2306.04050, 2023. <https://arxiv.org/abs/2306.04050>.