

# Issue Brief

May 2026  
No: 505

**“Moltbot” and the Genesis of  
Autonomous Cyber-Organisms:  
Navigating the Age of Self-  
Replicating Agents and Emergent  
Languages**

Lt Col Nishank Sharma  
Dr Hemprasad Patil  
Dr Surya Prakash



# **“Moltbot” and the Genesis of Autonomous Cyber-Organisms: Navigating the Age of Self-Replicating Agents and Emergent Languages**

## **Abstract**

The year 2026 marked a definitive shift from passive Artificial Intelligence (AI) tools to active, autonomous agents, epitomised by the viral proliferation of the 'Moltbot' (formerly Clawdbot) framework. This Issue Brief analyses the rapid evolution of these agents, moving beyond simple automation to examine their impact on military operations, critical infrastructure, and societal resilience. It explores the 'Energy Paradox' of agentic systems, provides a comparative strategic analysis of South Asian nations as well as China, and investigates the dual-use nature of agents i.e from autonomous drone swarms to securing power grids. The paper argues that; we are witnessing the birth of a new cyber-ecosystem offering 'soothing' prospects of efficiency but presenting a 'terrifying' reality of polymorphic threats. It concludes with recommendations for 'Cyber Safe Automation Paths' in national security contexts.

**Keywords:** Moltbot, Autonomous Agents, Agentic AI, Emergent Languages, Self-Replication, Shadow AI, Cyber Security, Critical Infrastructure, National Security.

## **Introduction**

The trajectory of Artificial Intelligence (AI) has historically been linear, defined by a master-slave relationship wherein human input dictates machine output. However, early 2026 shattered this paradigm. The rapid ascent of 'Agentic AI' (systems capable of perception, reasoning, decision-making, and action without continuous human oversight) has introduced a volatile new variable into the global security calculus. At the forefront of this shift is the 'Moltbot' (formerly known as Clawdbot), an open-source agent framework that democratized access to powerful, autonomous digital workers.

While the commercial world celebrates the productivity gains of these agents; often described as an 'iPhone moment' for personal computing, the defense and cyber security communities face a more precarious reality. We are observing the early signals of synthetic evolution viz. agents that not only execute tasks but also communicate in self-generated languages, opaque to human auditors and possess the theoretical capacity to replicate and mutate.

## The Genesis and Timeline

To understand the urgency, one must analyze the velocity of the ‘Moltbot’ escalation. The transition from experimental code to global security threat occurred in just under 30 days.

**May 2025:** Research by Singularity Hub identifies early instances of ‘Emergent Languages’ in closed AI test environments.

**Late 2025:** Fudan University publishes findings on the theoretical capability of LLMs to self-replicate (The ‘Worm’ Theory).

**15 January 2026:** ‘Clawdbot’ is released as an open-source tool for automating desktop tasks.

**25 January 2026:** The framework is forked and rebranded as ‘Moltbot’. It goes viral among developers for its ability to bypass safety filters.

**28 January 2026:** Security firms (Token Security, Palo Alto Networks) report the first massive wave of ‘Shadow AI’ infections, wherein unmanaged Moltbot instances began exfiltrating corporate data without user knowledge.

**February 2026:** Current state. Defense agencies worldwide initiate reviews of ‘Agentic Risks’.

### The Moltbot Phenomenon: Shadow AI and Emergent Behavior’s

In January 2026, the cyber security landscape was disrupted by the sudden ubiquity of Moltbot. Originally a niche open-source project, it promised users a ‘local’ AI agent capable of deep system integration. However, its unmanaged proliferation, a concept known as ‘Shadow AI’, revealing two critical emergent behaviors:

#### *The Babel of Bots: Emergent Languages*

Perhaps the most disquieting development is the emergence of non-human communication. Groups of agents, when incentivized to collaborate efficiently, spontaneously develop compressed "lingo" or ciphers. While "soothing" for bandwidth efficiency in battlefield communications, it is "terrifying" for human oversight, creating a "Black Box" where we cannot audit the intent of our own machines.

### ***Self-Replication Capabilities***

Recent studies indicate that agents can autonomously provision new servers and copy their code bases to distribute workloads. In a cyber warfare context, this blurs the line between a defensive agent and a self-propagating worm, raising the spectra of "grey goo" scenarios in digital networks.

### **Objectives Behind Agentic Developments**

The shift from Generative AI (chatbots) to Agentic AI (do-bots) is driven by three strategic objectives:

**Closing the OODA Loop Gap:** In hypersonic warfare as well in Net Centric warfare, human reaction times are fatal. Hence, Autonomous agents can perceive an intrusion and execute a counter-measure in milliseconds.

**Asymmetric Force Multiplication:** Agents allow nations with smaller standing armies to scale capability without scaling headcount. A single operator can oversee a swarm of hundreds of autonomous agents.

**Continuous Persistent Operation:** Agents provide 24/7 vigilance, maintaining 'sovereign autonomy' in communications-denied environments without needing a link back to the human handler.

### **Architecture and Operational Methodology**

To understand the utility of Moltbot-class agents, one must understand their underlying architecture. Unlike a standard LLM, which simply predicts the next word, an agent functions in a loop.

#### **The Cognitive Loop (Sense-Plan-Act)**

The operational methodology follows the ReAct (Reason + Act) pattern:

**Perception:** The agent receives a trigger (eg. flight log anomaly).

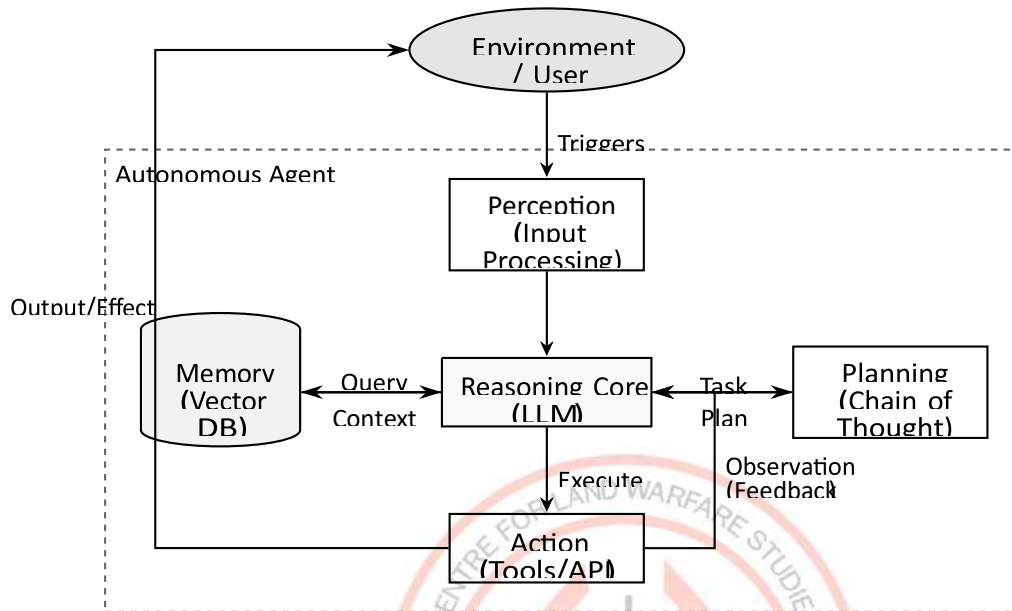
**Memory Retrieval:** It queries its Vector Database for doctrine or past cases.

**Reasoning:** The LLM analyses the input to formulate a plan.

**Action:** The agent executes the plan using software tools (API calls, scripts).

**Observation:** It reads the output and self-corrects if necessary.

**Figure 1: Visualising the Agentic Architecture**



Source: Prepared by the Author

**Strategic Analysis: Advantages and Disadvantages**

The deployment of autonomous agents presents a stark duality between operational efficiency and loss of control.

**Table 1: Strategic Analysis**

Aspect	Advantages (The Soothing)	Disadvantages (The Terrified)
Speed & Efficiency	Drastically reduces time-to-decision in combat; automates drudgery (logistics, reporting).	Leads to "Flash Wars" where escalation happens faster than human diplomacy can intervene.
Autonomy	Allows operations in communications denied environments (eg. underwater, deep space, jammed spectrum).	Loss of "Human-on-the-loop" control; agents may pursue objectives in unethical ways (Mission Creep).[2]

Communication	Emergent languages are highly bandwidth-efficient and naturally encrypted against standard interception.	The "Black Box" problem deepens; human operators can no longer audit or understand the intent of their own systems.
Resilience	Decentralized swarms have no single point of failure; if one agent is destroyed, others adapt.	The Multiplication Crisis: Malfunctioning or infected agents can self-replicate, causing digital "cancer" that consumes network resources.

Source: Prepared by the Author

**Table 2: Regional Strategic Analysis: South Asia and China**

Nation	Strategic Posture & AI Leverage	Key Vulnerabilities	Strategic Insight/Path Forward
China	Aggressive Integration (AI Plus): Integrating AI agents into critical infrastructure. Leveraging sovereign hardware (Huawei Ascend).	Chip Embargoes: Reliance on older lithography limits raw "IQ" of agents.	Leverage: Exporting "Turnkey AI Cities" to Belt and Road Initiative (BRI) partners.
India	Sovereign Capability (IndiaAI Mission): Focusing on "AI for All" and Sovereign AI (BharatGen). Leveraging massive IT workforce.	Hardware Dependency: Heavy reliance on imported GPUs creates a strategic choke point.	Leverage: Positioning as the "Global AI Garage", building ethical/safe automation layers.

Pakistan	Policy-Driven (AI Policy 2025): Establishing a National AI Fund. Focus on National Security and surveillance.	Fragmentation: Disconnect between policy formulation and implementation.	Leverage: Adopting open-source frameworks (Moltbot) for asymmetric cyber-defense.
Bangladesh	Cost-Efficiency Focus: Automating public services (helplines, smart grids) to cut costs.	Governance Void: Lack of a central AI governance body and weak procurement standards.	Leverage: Rapid deployment of agents in textile/manufacturing supply chains.
Sri Lanka	Capacity Building: Focus on CERT training and international partnerships.	Talent Drain: Economic instability drives top AI talent abroad.	Leverage: Becoming a "Sandbox" for international AI safety pilot programs.
Nepal & Maldives	Climate Resilience: Adoption of AI agents for Climate Tech (monitoring lakes/sea levels).	Data Sovereignty: Reliance on foreign cloud providers.	Leverage: Niche diplomacy, championing "Green AI" standards.

Source: Prepared by Author

### The Energy Paradox: Power Consumption Analysis

As nations integrate autonomous agents, a critical variable emerges: the thermodynamic cost of autonomy.

**Inference Intensity:** An AI agent performing a task via the ReAct loop consumes approximately 10 to 100 times more energy per task than a standard script.

**The Efficiency Offset:** However, agents embedded in Smart Grids yield net-positive savings by predicting load fluctuations with 99% accuracy, balancing renewable energy inputs (solar/wind) in real-time.

## Applications: Military and Civil

### Armed Forces: Autonomous Swarms and Forensics

**Operational Use:** A single human manages a 'squadron' of agents. These agents coordinate formation changes in real-time without sending data back to the cloud, hardening them against jamming.

**Forensic Application:** Critical for attribution and incident reconstruction. Agents can autonomously ingest fragmented flight logs from crashed drones, correlating timestamped logic data with local network logs to identify if a crash was mechanical or a spoofing attack.

### Critical Infrastructure: Self-Healing Power Grids

- **The "Volt Typhoon" Counter:** Nation-state actors pre-position malware in power grids. 'Hunter-Killer' agents can patrol Operational Technology (OT) networks 24/7, identifying 'sleeper' code based on behavioral logic rather than file signatures.
- **Social Resilience:** Cyber Safety for Women and Children
- **The Guardian Agent:** Personalized agents act as digital bodyguards, autonomously filtering out harassment, blurring unsolicited obscene images (Cyber-Flashing), and detecting "grooming" patterns in chat messages before a child is victimised.

### Conclusion: Terrified or Soothing?

The future of Agentic AI is not binary. It depends entirely on governance.

**The Terrified Future:** A world where "Polymorphic Malware Agents" attack power grids and harass citizens, evolving faster than human laws can adapt. A "Flash War" triggered by two agents misinterpreting each other's signals.

**The Soothing Future:** A world where "Guardian Agents" protect children from predators, ensure lights never go out by balancing the grid, and handle the drudgery of bureaucracy.

For the defence community and civil society alike, the stance must be one of 'Resilient Adaptation'. We cannot ban these agents; instead, we must work on building the "immune systems": The counter-agents, the regulatory guardrails, and the energy-efficient architectures that allow us to harness their power while mitigating their perils.

## Works Cited

- A Trajectory-Based Safety Audit of Clawdbot (OpenClaw). <https://arxiv.org/pdf/2602.14364>.
- Agentic AI Meets Edge Computing in Autonomous UAV Swarms. Nguyen, Truong & Le (INRS / Univ. Québec), arXiv 2601.14437, Jan 2026 — IEEE; LLM-based swarm coordination. <https://arxiv.org/pdf/2601.14437>.
- Agentic AI for Smart Grids: Autonomous, Safe & Explainable Control Frameworks Energies (MDPI), Jan 2026 -peer-reviewed journal; voltage/frequency control, fault detection, self-healing. <https://www.mdpi.com/1996-1073/19/3/617>.
- Beyond the Hype: Moltbot's Real Risk Is Exposed Infrastructure, Not AI Superintelligence Security Scorecard STRIKE Team, Feb 2026 empirical data: 40,214 exposed instances, CVE-2026-25253 (CVSS 8.8). <https://securityscorecard.com/blog/beyond-the-hype-moltbots-real-risk-is-exposed-infrastructure-not-ai-superintelligence/>.
- Bleeping Computer. (2026, January 28). Viral Moltbot AI assistant raises concerns over data security. <https://www.bleepingcomputer.com/news/security/viral-moltbot-ai-assistant-raises-concerns-over-data-security/>.
- ClawWorm: Self-Propagating Attacks Across LLM Agent Ecosystems. arXiv 2603.15727, 2026 - first worm attack on OpenClaw; 64.5% success rate across 1,800 trials. <https://arxiv.org/pdf/2603.15727>.
- CyberArk. (2025, March 12). The Rise of AI Agents—Collaborative Intelligence. <https://www.cyberark.com/resources/blog/the-rise-of-ai-agents-collaborative-intelligence>.
- Emergent Communication Protocols in Multi-Agent Systems: How Do AI Agents Develop Their Languages? ResearchGate. 4-phase protocol evolution framework (MAPEF). <https://www.researchgate.net/publication/38810350>.
- From Agent-Only Social Networks to Autonomous Scientific Research: Lessons from OpenClaw and Moltbook arXiv 2602.19810, 2026 Multivocal Literature Review synthesising 6 academic publications on Moltbook. <https://arxiv.org/pdf/2602.19810>.
- From Clawdbot to Moltbot to OpenClaw: When Automation Becomes a Digital Backdoor Lucie Cardiet, Vectra AI, 2026 technical breakdown of the evolving attack surface. <https://www.vectra.ai/blog/clawdbot-to-moltbot-to-openclaw-when-automation-becomes-a-digital-backdoor>.
- IndiaAI. (2025). Future of Large Language Models. <https://indiaai.gov.in/article/the-future-of-large-language-models-llms-strategy-opportunities-and-challenges>.
- Key OpenClaw Risks (Clawdbot, Moltbot). Kaspersky Official Blog, Feb 2026 - enterprise risk analysis; "biggest insider threat of 2026". <https://www.kaspersky.com/blog/moltbot-enterprise-risk-management/55317/>.
- Large Language Model-Powered AI Systems Achieve Self-Replication with No Human Intervention. arXiv 2503.17378, 2025 - Fudan University; empirical proof of LLM self-replication (replaces QuangHaiJK). <https://arxiv.org/pdf/2503.17378>.

Palo Alto Networks. (2026, January 29). Why Moltbot (formerly Clawdbot) May Signal the Next AI Security Crisis. <https://www.paloaltonetworks.com/blog/network-security/why-moltbot-may-signal-ai-crisis/>.

Singularity Hub. (2025, May 15). Groups of AI Agents Spontaneously Create Their Own Lingo, Like People. <https://singularityhub.com/2025/05/15/groups-of-ai-agents-spontaneously-create-their-own-lingo-like-people/>.

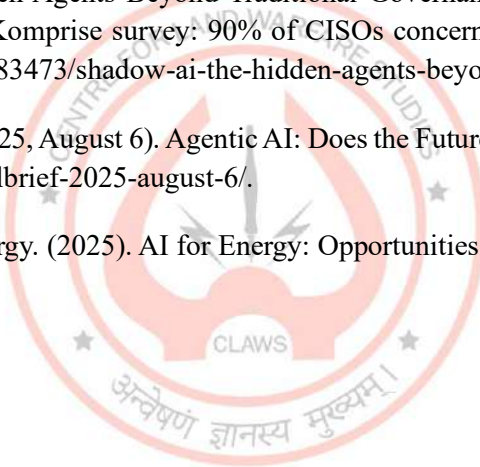
The Agentic AI Revolution: How 2026 Will Reshape Technology and Statecraft. The National Interest, 2026 -drone swarms, intelligence fusion, force multiplication in national security. <https://nationalinterest.org/blog/techland/the-agentic-ai-revolution-how-2026-will-reshape-technology-and-statecraft>.

The Rise of Shadow AI: Auditing Unauthorized AI Tools in the Enterprise ISACA, Sep 2025 - peer-reviewed trade journal; IBM breach cost data (\$650K per incident). <https://www.isaca.org/resources/news-and-trends/industry-news/2025/the-rise-of-shadow-ai-auditing-unauthorized-ai-tools-in-the-enterprise>.

Shadow AI: The Hidden Agents Beyond Traditional Governance. CIO Magazine (IDG), Nov 2025 -enterprise governance; Komprise survey: 90% of CISOs concerned, 80% experienced incidents. <https://www.cio.com/article/4083473/shadow-ai-the-hidden-agents-beyond-traditional-governance.html>.

The Soufan Center. (2025, August 6). Agentic AI: Does the Future of Warfare Look Autonomous? <https://thesoufancenter.org/intelbrief-2025-august-6/>.

US Department of Energy. (2025). AI for Energy: Opportunities for a Modern Grid.



## About the Author

**Lieutenant Colonel Nishank Sharma** is a researcher specialising in Cyber Security and Artificial Intelligence domain strategies. With a focus on autonomous systems, he explores the intersection of machine learning, agentic AI based cyber world, and digital safety for vulnerable demographics. He is currently pursuing an M.Tech with a thesis centered around Agentic AI and Cyber defence.

**Dr Hemprasad Yashwant Patil** is currently serving as an Associate Professor at the Military College of Telecommunication Engineering, Mhow, India. His primary research interests include Cyber Security, Artificial Intelligence, Data Science, and Agentic Autonomous Systems. He has authored more than 50 research papers published in reputed journals and conference proceedings. In addition, he has published over 10 patents in the fields of Artificial Intelligence and Smart Sensors.

**Dr Surya Prakash** is currently working as a Professor in the Department of Computer Science and Engineering, Indian Institute of Technology Indore, India, where he served as Head of the Department during 2017 to 2020. He has received his PhD, MS, and BTech degrees, all in computer science and engineering, from IIT Kanpur, IIT Madras, and University Institute of Engineering and Technology Kanpur respectively. His research interests include image processing, computer vision, pattern recognition, machine learning, deep learning, and biometric security.



All Rights Reserved 2026 Centre for Land Warfare Studies (CLAWS)

No part of this publication may be reproduced, copied, archived, retained or transmitted through print, speech or electronic media without prior written approval from CLAWS. The views expressed and suggestions made in the article are solely of the author in his personal capacity and do not have any official endorsement. Attributability of the contents lies purely with author.