


# The RSI Inflection Signal

What Anthropic's Internal Productivity Data Means for the Enterprise Threat Horizon

2026-06-15

 AI-assisted Rapid Research



**© 2026 Cloud Security Alliance. Some rights reserved.**

You may download, store, display, view, print, redistribute, and link to this document in its original, unmodified form, provided that attribution to the Cloud Security Alliance is maintained and all trademark and copyright notices remain intact.

This document may not be modified or altered. You may quote portions of the document as permitted by the Fair Use provisions of the United States Copyright Act, provided that attribution is given to the Cloud Security Alliance.

This document may be shared on professional and social media platforms in its original form with attribution.

*This document was generated with AI assistance and has not undergone official CSA review and approval processes.*

---

## Key Takeaways

- On June 4, 2026, Anthropic published "When AI builds itself," a detailed institutional report documenting that AI-assisted development has produced compounding productivity gains at Anthropic, and arguing that recursive self-improvement – AI systems improving their own successors – could arrive before most organizations are prepared to govern it [1].
- As of May 2026, more than 80% of code merged into Anthropic's production codebase was authored by Claude; by the second quarter of 2026, the median engineer was shipping 8× as much code per day as in 2024 – though Anthropic acknowledges this overstates true output quality, with an internal poll of 130 staff placing the genuine productivity multiplier at approximately 4× [1].
- The autonomous task horizon at Anthropic has been expanding rapidly: from 4-minute tasks in March 2024 to 90-minute tasks in March 2025 to 12-hour tasks in March 2026, with week-length tasks projected by 2027 [1]. The most recent interval – March 2025 to March 2026 – implies a doubling time of approximately four months, down from a prior rate of roughly seven months, though three data points are too few for a stable cadence claim.
- The same capability class demonstrated in Anthropic's defensive Project Glasswing – which surfaced more than 10,000 high- or critical-severity vulnerabilities across all partner-evaluated software in its first month, including 6,202 candidates across more than 1,000 open-source projects – is a dual-use capability that adversaries are likely to pursue through frontier model access, API misuse, or jailbreaks [2][3].
- Human code review has become the binding security constraint in AI-accelerated development: when generation velocity outpaces human review capacity, organizations lose their primary mechanism for detecting vulnerabilities, logic errors, and unintended behaviors introduced by AI-authored code [4].
- Anthropic's three-scenario framework – capability plateau, human-guided acceleration (assessed as most likely), and full recursive self-improvement – defines materially distinct risk profiles for enterprise threat modeling, and the human-guided acceleration scenario is already empirically underway [1].

# Background

## The "When AI Builds Itself" Report

On June 4, 2026, Anthropic's Institute published a document titled "When AI builds itself," with the subtitle "Our progress toward recursive self-improvement, and its implications" [1]. It is notable for explicitly disclosing internal productivity metrics, articulating the trajectory toward recursive self-improvement, and calling for international coordination to govern the outcome. The report draws on Anthropic's own engineering telemetry, internal surveys, and model evaluation benchmarks to make the case that AI-assisted AI development has crossed a qualitative threshold.

Recursive self-improvement, as the report defines it, describes a system capable of "making meaningful improvements to its own capabilities" in a compounding feedback loop, distinct from one-time human-directed training. Anthropic is explicit that this threshold has not yet been reached – current systems lack the independent research judgment and goal-selection required for genuine RSI – but the report's central argument is that the empirical trends put that threshold within a plausible near-term horizon, and that existing governance institutions are not positioned for the transition [1].

## The Internal Productivity Signal

The quantitative case in the report rests on Anthropic's own engineering operations. Before Claude Code launched in research preview in February 2025, Claude authored code in the low single digits as a share of merged production commits. By May 2026, that figure had risen to more than 80% [1]. The report frames this as an inflection rather than a gradual drift: productivity metrics held roughly constant across Anthropic's first four years of operation (2021–2024), then began to climb sharply as Claude transitioned from suggestion to execution, and steepened further when models began operating autonomously over longer time horizons.

The 8× code-per-engineer figure is presented with an explicit caveat. Lines of code measures quantity, not quality, and Anthropic's internal March 2026 survey – polling 130 employees across research teams – placed the median self-reported output multiplier at approximately 4× when using Mythos Preview compared to working without AI assistance [1]. That figure is nonetheless significant: a 4× sustained productivity multiplier across engineering and research functions implies a compression of the capability development calendar that may not yet be reflected in most enterprise risk model refresh cycles.

## The Task Duration Trajectory

One of the most operationally significant data series in the report tracks the maximum reliable autonomous task duration for successive Claude models. In March 2024, Claude Opus 3 could complete tasks requiring roughly four minutes of human work. By March 2025, Claude Sonnet 3.7 had extended that horizon to 90-minute tasks. By March 2026, Claude Opus 4.6 handled 12-hour tasks with comparable reliability [1]. Evaluations using the METR task complexity benchmark recorded Claude Mythos Preview sustaining autonomous work for at least 16 hours, approaching the upper end of what METR can measure without additional instrumentation [1].

The most recent interval – March 2025 to March 2026 – implies a doubling time of approximately four months, down from a prior rate of roughly seven months, meaning the doubling cadence itself has been accelerating [1]. Anthropic does not assert that this rate will hold, and three measured data points are a thin basis for claiming a stable cadence. Nonetheless, even a moderated version of this trajectory has direct implications for the threat planning horizons that enterprise security teams apply to agentic AI risk. If the four-month doubling held through 2026, reliable multi-day autonomous task completion would arrive before year-end; week-length tasks would follow in 2027.

## Security Analysis

### Threat Velocity Compression

The enterprise security implication most immediately derivable from Anthropic's data is what might be called threat velocity compression: as AI systems become capable of completing longer and more complex tasks autonomously, the time required to develop and deploy an exploit shrinks in proportion. Capabilities that once required a skilled human team working for days – reverse engineering a binary, constructing a working exploit chain, identifying lateral movement paths across a network – increasingly fall within the autonomous task horizon of frontier models.

Anthropic's own Project Glasswing demonstrates the ceiling of this capability in a controlled defensive context. The program granted a small set of vetted defensive partners exclusive access to Claude Mythos Preview; within its first month of operation, it surfaced 6,202 candidate vulnerabilities across more than 1,000 open-source projects, assessed 1,752 of those, and validated 1,587 as true positives, of which 1,094 were confirmed as high- or critical-severity [2][3]. Among the confirmed findings was CVE-2026-5194, a CVSS 9.1 flaw in the WolfSSL library that enables certificate forgery and service impersonation [2]. The detection was autonomous: no human researcher directed Mythos to that specific codebase or flaw class.

Glasswing establishes what an authorized program with frontier model access can achieve on defense. The same underlying capability – autonomous vulnerability discovery at scale – is dual-use. While no public incident has confirmed an adversary deploying equivalent autonomous offensive capability, the capability demonstrated by Glasswing is not exclusive to defenders. Nation-state threat actors and well-resourced criminal organizations have attempted to acquire frontier AI capabilities through jailbreaking, credential theft, and illicit API access. An organization that treats its current-generation threat model as stable across a 12-month planning horizon is implicitly assuming that adversary AI capability will not advance materially in that interval. Anthropic's data suggests that assumption has become difficult to defend.

## The Code Review Bottleneck as Security Gap

In AI-accelerated development environments, human code review has emerged as the rate-limiting control in the security architecture – and one whose capacity has not scaled proportionally with AI generation velocity. When 80% of production code is AI-authored and a single engineer is merging 8× the volume they managed in 2024, the review capacity that traditional security programs sized for human development throughput is no longer adequate [1]. Bugs, logic errors, and adversarially introduced behaviors that a patient reviewer would catch may pass undetected simply because no proportional review capacity exists to surface them.

This creates a compounding vulnerability: AI-generated code at high volume introduces the statistical likelihood of undetected defects, while simultaneously reducing the time available per review cycle. Organizations that have deployed AI coding assistants without correspondingly expanding automated code analysis, semantic security testing, or AI-assisted review tooling have structurally weakened a previously reliable security control without recognizing that they have done so [4].

Anthropic's report is direct about this within its own operations. The company notes that human approval of code and experiments "is the mechanism by which organizations stay accountable," and describes human review as a new constraint – not a residual overhead – in the face of machine-speed generation [1]. The appropriate enterprise response is not to slow AI adoption, but to recognize that code review tooling must be upgraded alongside it.

## The Recursive Third-Party Risk Problem

Anthropic's RSI report introduces a risk category that conventional third-party due diligence frameworks are not equipped to assess. When a vendor uses its own AI system to develop that same AI system's successors, the standard security evaluation questions – "what is their engineering team's credential?"

"what is their code review process?" "how do they manage SDLC security?" – become category errors. The answers are partially or substantially determined by the behavior of the product being sold [4].

For enterprise organizations that have integrated Claude Code, Claude API, or related Anthropic tooling into development pipelines, this creates an auditability challenge without a clear precedent. A defect or behavioral anomaly in the current generation of Claude models may influence the successor model's behavior via the training pipeline. Code quality issues, security boundary misunderstandings, or subtle misalignment in Anthropic's current codebase – 80%+ of which was AI-authored – could propagate into future model versions. Traditional audit methodologies lack the instrumentation to detect or characterize this risk class.

Anthropic's report explicitly acknowledges this uncertainty [1]. As Bellamkonda observes in his analysis of the report, the core concern is what happens "when the system building the next model also has bugs" [4]. This acknowledgment does not resolve the practical enterprise risk management question of how to account for such recursive risk in vendor assessments and dependency reviews.

## **Asymmetric Adoption and the Defender Gap**

Enterprise security organizations face a structural disadvantage in the current AI acceleration environment. Frontier model capabilities, by definition, are available first to frontier labs and their immediate partners, then diffuse to enterprise adopters over a period of months to years. The capability that Glasswing partners deployed for defensive vulnerability discovery in May 2026 reflects a model generation that adversaries are simultaneously attempting to access for offensive use – and the defensive deployment required selective, vetted access through a controlled program, while offensive use requires a capable jailbreak, stolen credential, or shadow API access – each lower-friction than operating an authorized institutional program like Glasswing.

The asymmetry is not inevitable. Organizations that invest in AI security tooling during this window – before the doubling cadence places the next capability level within adversaries' reach – accrue a compounding defensive advantage. Those that defer AI security investment on the grounds that current-generation AI threats are manageable with existing controls are implicitly betting against the observed trajectory. The Forrester 2026 Threat Intelligence Report identified AI agents as the top CISO risk category for the year, with near-autonomous nation-state attacks and AI identity sprawl among the most pressing operational concerns (Forrester 2026 Threat Intelligence Report, as reported by Cybersecurity Insiders [5]).

## Scenario Planning and Planning Horizon Calibration

Anthropic's three-scenario framework provides a structured basis for enterprise security planning that is more useful than binary predictions about AI timelines [1][8]. In Scenario 1, the capability plateau, the current productivity metrics represent a ceiling rather than a trajectory; enterprise security teams should plan for high-volume AI-generated code and faster exploit tooling, but the threat model stabilizes within the current capability band. In Scenario 2, human-guided acceleration – which Anthropic assesses as the most likely near-term path – AI systems continue to assume a growing share of software development and research while humans retain directional control; the threat model must accommodate a recurring capability step-up on an accelerating cadence, requiring correspondingly frequent threat model revision cycles [1].

Scenario 3, full recursive self-improvement, describes a qualitatively different risk environment in which AI systems design their own successors with minimal human involvement and the pace of development is bounded primarily by compute availability. Anthropic does not assess this scenario as imminent, but it is explicit that it cannot be ruled out, and it identifies the remaining gap – independent research judgment and goal selection – as a capability threshold rather than a categorical barrier [1]. Enterprise security teams should not treat Scenario 3 as the baseline for current planning, but they should have a position on what governance posture they would adopt if evidence of that transition emerged.

## Recommendations

### Immediate Actions

Security teams should inventory all AI-assisted code generation currently operating in their development environment and audit the code review controls applied to AI-authored commits. Organizations should treat the gap between AI generation velocity and review capacity as an unmitigated control exposure unless they have verified that complementary controls – automated analysis, AI-assisted review, policy restrictions – are fully compensating. Static analysis, dynamic testing, and AI-assisted review tooling should be assessed and deployed to close that gap before it expands further as developer adoption grows.

Organizations should also conduct a targeted vendor risk reassessment for any AI provider whose model is used in development pipelines. The standard third-party security questionnaire is insufficient for AI vendors engaged in AI-assisted self-development. Supplementary questions should address: the share

of production code authored by AI, the controls applied to AI-authored code before production merge, the auditability methodology for AI-influenced model training, and the incident response procedures specific to AI-generated behavioral anomalies.

## Short-Term Mitigations

Enterprise threat model refresh cycles should be shortened from annual or biannual to quarterly, with explicit checkpoints tied to frontier model generation releases. A threat model that does not account for the accelerating doubling cadence in autonomous task capability will systematically underestimate attacker capability by the time it informs operational security decisions. Each revision should assess whether existing detection and response controls remain adequate against the most recently confirmed autonomous task horizon.

AI code provenance tracking – the capacity to determine what portion of a production system was authored, modified, or reviewed by AI – should be established as a standard engineering hygiene practice and a security audit input. Without provenance data, organizations cannot assess the exposure created by AI code review gaps, cannot characterize incident root causes that may involve AI-authored code, and cannot comply with emerging AI transparency requirements in regulated sectors.

Organizations should also establish an AI red team exercise program with explicit scope covering autonomous agent attack chains and AI-assisted vulnerability discovery. The Glasswing findings – 1,094 confirmed high- or critical-severity vulnerabilities discovered autonomously across open-source projects in one month [2][3][7][10] – suggest that internal assets are likely to contain undiscovered flaws at comparable density. Organizations that wait for adversaries to surface these flaws forfeit the remediation window that a proactive AI-assisted internal assessment would provide.

## Strategic Considerations

At the strategic level, enterprise security programs should develop explicit positions on each of Anthropic's three RSI scenarios and define the governance posture and capability investments appropriate to each. Scenario 2, the most likely near-term path, implies that the current enterprise AI security investment cycle is not a one-time adjustment to AI adoption but an ongoing refresh program calibrated to an accelerating capability cadence. Budgeting and staffing models that treat AI security as a project rather than a continuous function will fall structurally behind the threat trajectory.

Security leadership should also engage with AI capability monitoring as an emerging intelligence function. The METR autonomous task benchmark, AI coding performance benchmarks such as SWE-bench, and vendor-published productivity data provide leading indicators of the threat capability

environment approximately six to twelve months before that capability diffuses to adversary toolsets. Organizations that build monitoring into their intelligence programs gain planning lead time that reactive approaches forfeit.

Finally, enterprises should follow and engage with the international governance mechanisms that Anthropic's report calls for, including multi-stakeholder deliberation processes and verification systems for AI development coordination [1][9]. Anthropic's stated willingness to implement a verifiable development slowdown contingent on reciprocal commitments from other frontier labs represents a governance approach that, if realized, would directly shape the pace of capability diffusion into the threat landscape. Enterprise security teams have standing interests in that outcome and should make those interests visible in industry and government consultation processes.

## CSA Resource Alignment

Anthropic's RSI data is most directly relevant to threat modeling work under CSA's MAESTRO framework. MAESTRO's seven-layer architecture addresses AI threats from the foundation model layer through ecosystem integration, and the RSI inflection signal touches at least three of those layers concretely [6]. Layer 1 (Foundation Model) encompasses the risk that model-level defects or misalignments propagate into downstream applications – a risk now complicated by the recursive feedback loop in which Claude authors a growing share of Anthropic's own code and training pipeline tooling. Layer 3 (Agent Frameworks) covers the reasoning loops and tool dispatch mechanisms that govern autonomous task execution, and the accelerating task horizon doubling directly recalibrates the risk surface that Layer 3 controls must address. Layer 7 (Ecosystem Integration) encompasses supply chain and third-party dependencies, which the recursive vendor risk problem implicates in a novel way.

CSA's AI Controls Matrix (AICM) provides a control mapping framework for managing AI-specific risks that can be applied to both the code review gap and the AI provenance tracking recommendations above. AICM, as a superset of the Cloud Controls Matrix, addresses data governance, model transparency, and operational accountability requirements that align directly with the audit and review controls described in this note's recommendations.

The AI Organizational Responsibilities publications from CSA – including the governance, risk management, and compliance guidance – provide the organizational accountability structures within which the strategic recommendations above should be implemented. The core principle that human approval remains the accountability mechanism for AI-generated output, which Anthropic's report articulates directly, is operationalized through the governance structures that those publications address.

Enterprises building AI security programs responsive to the RSI inflection signal should also engage with CSA's STAR for AI initiative, which provides a structured framework for communicating AI capability and risk posture to customers, partners, and regulators – and which creates a basis for the kind of reciprocal accountability that meaningful industry coordination requires.

## References

- [1] Anthropic. ["When AI builds itself: Our progress toward recursive self-improvement, and its implications."](#) Anthropic Institute, June 4, 2026.
- [2] Anthropic. ["Project Glasswing: An initial update."](#) Anthropic Research, May 2026.
- [3] Kovacs, Eduard. ["Project Glasswing has uncovered 10,000 vulnerabilities: Anthropic."](#) CSO Online, May 2026.
- [4] Bellamkonda, Shashi. ["Anthropic's RSI Report Is the Most Honest Vendor Document in Enterprise AI. Read It Carefully."](#) shashi.co, June 2026.
- [5] Cybersecurity Insiders. ["Forrester 2026 Threat Intelligence Report: AI Agents Top CISO Risk List."](#) Cybersecurity Insiders, 2026.
- [6] Cloud Security Alliance. ["Agentic AI Threat Modeling Framework: MAESTRO."](#) CSA Blog, February 6, 2025.
- [7] The Hacker News. ["Claude Mythos AI Finds 10000 High-Severity Flaws."](#) The Hacker News, May 2026.
- [8] MindStudio. ["Anthropic RSI Report: Three Scenarios for the Future of AI and What They Mean for Builders."](#) MindStudio Blog, June 2026.
- [9] Shumaker, Loop & Kendrick. ["Client Alert: Anthropic's Call for a Global AI Pause: What Businesses Need to Know About the Governance Landscape."](#) Shumaker Law, 2026.
- [10] SecurityAffairs. ["Anthropic's Glasswing: 10000+ Vulnerabilities Found in One Month and the Patching Problem Has Never Been More Obvious."](#) SecurityAffairs, May 2026.