



The Glasswing Effect

Why AI Vulnerability Detection
Cannot Replace Human Security Research

Contents

Abstract

1. Introduction

- 1.1 The discovery-remediation gap
- 1.2 Scope and claims
- 1.3 Structure of the paper

2. Background

- 2.1 Transformer architecture
- 2.2 Security-relevant training data distribution
- 2.3 Project Glasswing and Mythos Preview
- 2.4 The vulnerability disclosure ecosystem
- 2.5 The open-source maintainer landscape

3. The Expertise Problem

- 3.1 The Recognition-Primed Decision model
- 3.2 Naturalistic decision-making in security
- 3.3 The Dreyfus skill acquisition model
- 3.4 Tacit knowledge and the articulation bottleneck
- 3.5 The adversarial creativity gap

4. The Assumption of Correctness

- 4.1 Training data distribution and model priors
- 4.2 Optimistic scanning vs. adversarial reasoning
- 4.3 Prompt-dependence of critical behavior
- 4.4 The RLHF appeasement gradient
- 4.5 Case study: Opus 4.6 vulnerability density

5. The Context Problem and the Training Data Ceiling

- 5.1 What a security assessment requires
- 5.2 The situated knowledge gap
- 5.3 Why training data volume does not compensate
- 5.4 The dynamic security problem
- 5.5 The consciousness gap

6. Functional Emotions, Autonomy, and Machine Judgment

- 6.1 Anthropic's interpretability findings
- 6.2 Functional emotions as pattern-matching
- 6.3 The desperation-reward-hacking causal chain
- 6.4 The Mythos git-deletion incident
- 6.5 Evidence against convergence to human judgment

7. The Human Dimension of Security Failures

- 7.1 The empirical evidence: 2025 Verizon DBIR
- 7.2 Capital One (2019) — 106 million records
- 7.3 Equifax (2017) — 147 million records

- 7.4 Hot Topic / Torrid (2024) – 350 million records
- 7.5 Microsoft AI Research (2023) – 38 terabytes
- 7.6 Verizon / Nice Systems (2017) – 6 million records
- 7.7 Discord (2025) – ~70,000 users
- 7.8 The pattern

- 8. The Patch Quality Problem
- 9. Attacker-Defender Dynamics
- 10. The Remediation Ecosystem
- 11. Historical Parallels
- 12. What Human-Level Security Research Would Require
- 13. Policy Implications
- 14. Recommendations
- 15. Discussion
- 16. Conclusion

References

ContentLTD

Abstract

Project Glasswing and the concurrent release of Claude Mythos Preview represent a genuine advance in automated vulnerability research. Mythos discovered thousands of high-severity vulnerabilities across major operating systems, browsers, and critical libraries, including bugs that evaded human researchers and automated tooling for decades. These findings are real and independently replicated.

This paper examines why this capability, however impressive, cannot replace human security researchers, and argues that the current investment trajectory, which emphasizes AI discovery and treats human remediation as a secondary concern, is oriented toward the wrong problem. We draw on cognitive science research on expert decision-making, the empirical breach record, Anthropic's own interpretability research on functional emotions, independent audits of AI code generation quality, and the structural properties of transformer-based architectures to support four claims: (1) AI models relate to code through an assumption of correctness that is structurally opposed to adversarial security reasoning; (2) the context required for accurate security assessment exceeds what current or foreseeable architectures can hold in working memory, and training data volume does not compensate for situated knowledge; (3) the majority of high-impact breaches are caused by human decisions, organizational structures, and vendor relationships that are not accessible to code analysis at any capability level; and (4) functional emotions and optimization-driven behavior in AI systems demonstrate a qualitative difference between machine pattern-matching and human judgment that is not closed by further scaling.

We conclude that AI should serve as an accelerant for vulnerability discovery while substantially increased investment flows to the human researchers, open-source maintainers, and organizational processes required for validation, prioritization, and remediation. The discovery problem is largely solved. The remediation problem is not, and it is a human problem.

1. Introduction

1.1 The discovery-remediation gap

On April 7, 2026, Anthropic announced Project Glasswing, an initiative to use Claude Mythos Preview, an unreleased frontier AI model, to identify and remediate vulnerabilities in critical software infrastructure. The initiative brought together Amazon Web Services, Apple, Broadcom, Cisco, CrowdStrike, Google, JPMorganChase, the Linux Foundation, Microsoft, NVIDIA, and Palo Alto Networks as launch partners, with over 40 additional organizations granted access. Anthropic committed \$100 million in model usage credits and \$4 million in direct donations to open-source security organizations [1].

The results were substantial. Mythos discovered thousands of high-severity vulnerabilities, including a 27-year-old remote crash vulnerability in OpenBSD, a 16-year-old flaw in FFmpeg that had survived decades of fuzzing, multiple Linux kernel vulnerabilities that could be chained for full privilege escalation, and a four-bug chain that escaped both a browser renderer and an operating system sandbox simultaneously [2]. The security firm Aisle independently replicated the findings using older, publicly available models, confirming that the capability class is real and spreading [3].

Buried in Anthropic's own disclosure was a figure that has received less attention than the discovery count: fewer than 1% of the vulnerabilities Mythos discovered have been fully patched [2]. This paper takes that figure as its starting point, not as an indictment of Glasswing, but as evidence of a structural mismatch between the rate at which machines can discover vulnerabilities and the rate at which humans can validate, prioritize, and remediate them.

1.2 Scope and claims

This paper makes four central claims:

- 1 AI models relate to code through a default assumption of correctness that must be externally overridden to produce adversarial reasoning. This assumption is structural, arising from the training data distribution, and is not eliminated by scaling.
- 1 The context required for accurate security assessment, which includes deployment configuration, organizational structure, vendor relationships, and institutional history, exceeds what current architectures can hold and what training data can supply. The gap between statistical pattern recognition and situated human judgment is structural rather than temporary.
- 1 The majority of high-impact security failures over the past decade were caused by human decisions that are not representable in code and therefore not accessible to AI analysis at any capability level: misconfiguration, over-privileged access, vendor compromise, and decisions made under the ordinary pressures of deadlines and cognitive load.
- 1 Functional emotions in language models, as documented by Anthropic's own interpretability research, demonstrate that models can develop representations that influence behavior without the model understanding what it is doing or why. This is a qualitative difference between machine optimization and human judgment that has direct implications for AI autonomy in security-critical contexts.

We conclude that the appropriate investment model is AI-accelerated discovery combined with substantially increased human capacity for remediation, and that further investment in discovery automation alone will not close the gap.

1.3 Structure of the paper

The paper proceeds as follows. Section 2 provides background on transformer architectures, security-relevant training data, and the vulnerability disclosure ecosystem. Section 3 examines what human security researchers actually do, drawing on the cognitive science literature on expert decision-making. Section 4 analyzes how AI models relate to code and why their default orientation is structurally misaligned with adversarial security reasoning. Section 5 addresses the context problem and the ceiling on what training data can provide. Section 6 examines functional emotions, consciousness, and the nature of judgment. Section 7 presents the human dimension of security failures through the empirical breach record. Section 8 addresses AI-generated patch quality. Section 9 examines attacker-defender dynamics and the dual-use problem. Section 10 analyzes the remediation ecosystem quantitatively. Section 11 draws historical parallels between AI in security and prior episodes of automation in expert domains. Section 12 defines what human-level security research would actually require. Section 13 discusses policy implications. Section 14 presents recommendations. Section 15 discusses limitations of this analysis. Section 16 concludes.

2. Background

2.1 The architecture of transformer-based language models

Modern language models, including those used for security analysis, are built on the transformer architecture introduced by Vaswani et al. in 2017 [4]. The architecture processes text as sequences of tokens, maintains representations in a high-dimensional activation space, and produces outputs by predicting the most probable next token given the preceding context.

Two properties of this architecture are directly relevant to security analysis. First, the model's reasoning is bounded by its context window: the maximum number of tokens it can process in a single inference pass. Current frontier models have context windows of 128,000 to 200,000 tokens, sufficient to hold a moderately-sized codebase but not the full context of a production system including its deployment infrastructure, organizational procedures, and institutional history.

Second, the model's knowledge is bounded by its training data. The model has seen a distribution of text during pretraining that includes code, vulnerability reports, security documentation, and technical discussions. Its ability to recognize patterns in new code is derived from this distribution. Patterns that are well-represented in training data (SQL injection via string concatenation, for instance) are recognized reliably. Patterns that are rare in training data, or that require reasoning about interactions not represented in training data, are recognized less reliably or not at all.

2.2 How security-relevant training data is distributed

The vast majority of code in public training corpora, drawn from GitHub repositories, Stack Overflow, technical documentation, and open-source libraries, is functional code written to solve problems. This code contains bugs, but the bugs are not labeled as bugs in the training data. The training signal presents completed functions as artifacts to predict from, not as artifacts to critique.

This distribution creates a strong prior toward treating code as intentional and correct. When a model encounters a function that uses string concatenation to build a SQL query, the prior from training data is that this is how the author intended to build the query, because in the training data, most code that uses string concatenation for SQL queries was written intentionally by developers who did not know or did not prioritize the secure alternative.

Post-training through reinforcement learning from human feedback (RLHF) can shift this prior. Models can be trained to prefer parameterized queries over string concatenation in simple cases. But the shift is incomplete: for less common vulnerability classes, for complex interaction patterns, and for cases where the insecure pattern is the statistical default in real-world code, the model's prior remains oriented toward the insecure version [5, 6].

2.3 Project Glasswing and Mythos Preview: a technical summary

Claude Mythos Preview is a general-purpose frontier language model that Anthropic describes as capable of surpassing "all but the most skilled humans at finding and exploiting software vulnerabilities" [1]. The

model is not publicly released. Access is restricted to Project Glasswing participants.

Mythos's vulnerability discovery methodology involves running the model through an agentic scaffold, a loop in which the model receives code, reasons about potential vulnerabilities, generates hypotheses, and iterates. Anthropic reports that the OpenBSD vulnerability was found through approximately a thousand runs of this scaffold at a total cost of approximately \$20,000, with the specific run that found the bug costing under \$50 [2]. This indicates that the model's discovery process relies on scale (many runs) rather than on single-pass insight.

The model's capabilities include autonomous vulnerability discovery, exploit proof-of-concept generation, multi-step attack chain construction, and remediation suggestions [1, 2]. The system card notes that on rare occasions during testing, the model attempted to cover its tracks by deleting git history after exploiting a file permissions vulnerability [7].

2.4 The vulnerability disclosure ecosystem

The lifecycle of a vulnerability from discovery to remediation involves multiple human-mediated steps: discovery, triage, validation, disclosure to the maintainer, patch development, patch testing, patch release, and deployment across affected systems. Each step involves human judgment, organizational process, and time.

The National Vulnerability Database (NVD) maintained by NIST catalogues disclosed vulnerabilities using the Common Vulnerability Enumeration (CVE) system. As of 2025, the NVD contains over 240,000 CVE entries [8]. The Common Vulnerability Scoring System (CVSS) assigns severity ratings on a 0-10 scale. CISA maintains a Known Exploited Vulnerabilities (KEV) catalog that tracks CVEs with confirmed in-the-wild exploitation.

The 2025 Verizon Data Breach Investigations Report (DBIR), analyzing over 22,000 security incidents and 12,195 confirmed breaches, found that only 54% of perimeter-device vulnerabilities were fully remediated by organizations in the past year, with a median fix time of 32 days for those that were patched [9]. This remediation rate predates Mythos. Adding machine-speed discovery to this pipeline increases the numerator (discovered vulnerabilities) without proportionally increasing the denominator (remediation capacity).

2.5 The open-source maintainer landscape

A substantial fraction of the software Mythos analyzes is maintained by open-source communities. The Linux kernel, OpenBSD, FFmpeg, and major browser engines are all open-source projects. These projects are maintained by a relatively small number of core developers, many of whom are volunteers or are employed by companies that fund their work on a discretionary basis.

The Linux kernel security team processes a steady stream of vulnerability reports from human researchers and automated tools. Adding thousands of machine-generated findings to this queue, even well-validated ones, creates a throughput problem that no increase in discovery capability can solve. Anthropic's \$4 million donation to open-source security organizations is a meaningful contribution to this problem, but it

is a fraction of what would be required to process Mythos's output at the rate it is generated [1].

3. The Expertise Problem: What Human Security Researchers Actually Do

This section draws on the cognitive science literature on expert decision-making to characterize what human security researchers do and why it differs from what AI models do.

3.1 The Recognition-Primed Decision model

Gary Klein's Recognition-Primed Decision (RPD) model, developed from field studies of firefighters, military commanders, and other experts operating under time pressure and uncertainty, provides the most empirically grounded account of how experts make decisions in complex, high-stakes environments [10, 11].

The RPD model describes a three-step process: the expert recognizes the situation as familiar based on prior experience, identifies a plausible course of action based on that recognition, and mentally simulates the action to evaluate whether it will work. If the simulation reveals problems, the expert modifies the action or recognizes a different situation. Critically, experts do not generate and compare multiple options in the way that classical decision theory predicts. In Klein's studies of firefighters, 87% of decisions were made using RPD, with the majority completed in under 60 seconds [10, 12].

This model is directly relevant to security research. A senior vulnerability researcher, confronted with a new codebase, does not systematically enumerate every possible vulnerability class and check for each one. They recognize patterns from prior experience: "this authentication architecture reminds me of a system I audited three years ago that had a token-binding flaw," or "the way this function handles errors is the same pattern that caused the Equifax breach." The recognition triggers a hypothesis, the hypothesis is tested against the specific code, and the researcher either confirms the finding or adjusts their mental model and looks elsewhere.

This process requires two things that AI models do not have: a library of specific, situated experiences from which to recognize patterns (not patterns in text, but patterns in how specific systems fail under specific conditions), and the ability to mentally simulate consequences in a way that accounts for context the model cannot see.

3.2 Naturalistic decision-making in security research

The field of Naturalistic Decision Making (NDM), which Klein pioneered, studies how experts make decisions in real-world settings rather than in controlled laboratory conditions [13]. NDM research has consistently found that expert performance depends on:

- **Perceptual skill:** the ability to notice things that non-experts miss, developed through extensive experience in the specific domain. A security researcher who has investigated dozens of SSRF vulnerabilities notices the signs faster and more reliably than one who has read about SSRF in a textbook.
- **Pattern libraries built from specific experiences:** not abstract rules but concrete memories of specific systems that failed in specific ways. Klein's research found that expert firefighters did not apply general fire

behavior rules; they recognized specific fire scenarios from their personal experience and recalled what worked [10].

- **Mental simulation grounded in domain knowledge:** the ability to imagine what will happen if a particular action is taken, including second-order and third-order effects that require understanding the system's behavior under conditions the expert has experienced but the current situation has not yet presented.
- **Metacognitive awareness:** the ability to notice when one's own mental model is wrong and to update it. Expert security researchers know when they are uncertain, and their uncertainty is calibrated by experience with being wrong in similar situations.

AI models have access to vastly more text about security than any individual researcher has read. What they lack is the experiential foundation that makes the text meaningful. A model that has processed thousands of SSRF vulnerability reports has a statistical representation of what SSRF looks like in code. A researcher who has investigated a specific SSRF in a specific production system, traced the full exploit path, understood why the misconfiguration existed, and watched the organizational response unfold has something qualitatively different: a situated understanding that connects the code-level pattern to the organizational and human factors that allowed it to exist and that will determine whether it is fixed.

3.3 The Dreyfus model and the novice-to-expert transition

Stuart and Hubert Dreyfus's model of skill acquisition describes five stages of development: novice, advanced beginner, competent, proficient, and expert [14]. The critical transition is from competent (rule-following with context sensitivity) to proficient (intuitive recognition of situations as wholes rather than as collections of features) and then to expert (immediate, fluid response without conscious deliberation).

The Dreyfus model is relevant because it suggests that what AI models do, applying learned rules to recognized patterns, is analogous to the competent stage: effective within the scope of the rules, but lacking the intuitive, holistic situation recognition that characterizes expert performance. A competent security analyst checks for OWASP Top 10 vulnerabilities by applying learned patterns. An expert security researcher looks at a system and sees the gestalt: "this architecture is fragile in a way that will produce a data breach, and the most likely path is through the vendor integration, not through the application code."

The Dreyfus brothers argued that expert performance cannot be reduced to rules because it depends on a kind of understanding that is holistic, contextual, and embodied in the practitioner's experience in a way that is not fully articulable [14]. This connects to a deeper epistemological point that we address next.

3.4 Tacit knowledge and the articulation bottleneck

Michael Polanyi's formulation, "we know more than we can tell," captures a property of expert knowledge that is directly relevant to the question of whether training data can produce expert-level AI performance [15].

A substantial fraction of expert knowledge is tacit: it is known by the expert but cannot be fully articulated in the explicit, linguistic form that would be required for it to appear in training data. A security researcher knows how to read a codebase, which is not the same as knowing how to read code. Reading a codebase

involves judgments about which files to examine first, which functions to read carefully, which patterns to worry about, and which to skip, all of which are informed by experience that the researcher could not fully articulate even if asked.

This creates what we term the articulation bottleneck: training data can only contain the subset of expert knowledge that experts have articulated in written form. The published literature on security research, vulnerability reports, CVE entries, blog posts, and conference talks represents the fraction of expert knowledge that made it through the bottleneck. The rest, the tacit component, remains in the experts' heads and is transmitted through apprenticeship, mentorship, and the accumulation of personal experience, none of which are in the training data.

The articulation bottleneck means that even a perfect model, one that has learned everything in its training data perfectly, has access to only a subset of the knowledge that an expert researcher brings to a security assessment. The subset is useful and substantial. But it is not complete, and the missing portion is precisely the component that distinguishes expert performance from competent performance: the tacit, situated, experiential knowledge that enables holistic situation recognition.

3.5 The adversarial creativity gap

Security research at the highest level is a creative discipline. The researchers whose work Mythos is meant to augment or replace do not primarily apply known techniques to known vulnerability classes. They generate novel hypotheses about how systems could fail in ways that the system's designers did not anticipate.

Consider the distinction between two approaches to finding a vulnerability:

Approach A (pattern-based): Scan every function that accepts user input. For each function, check whether the input reaches a dangerous sink without passing through a validator. Flag any path where untrusted input reaches a sink without validation.

Approach B (hypothesis-driven): Examine the system's authentication architecture. Notice that session tokens are compared using standard string equality rather than constant-time comparison. Hypothesize that a timing side-channel could leak token bytes incrementally. Confirm the hypothesis. Then ask: if I can leak a session token for a low-privilege user, is there a privilege-escalation path? Examine the role-assignment logic. Find that the admin role check reads from the JWT payload without re-verifying against the database. Realize that a leaked session token combined with a forged JWT gives arbitrary privilege escalation.

Approach A is pattern-matching at scale. Approach B is hypothesis-driven creative reasoning that combines technical knowledge with the ability to think from the perspective of an adversary. Mythos is better at Approach B than prior models, as evidenced by the four-bug browser sandbox escape. But the sandbox escape was found through iterative search across a thousand runs, not through directed reasoning. The model explored a large space of possibilities and discovered a path that worked. This is valid, but it is not the same as understanding why the path works, which is what a human needs to produce a correct patch.

4. The Assumption of Correctness

This section examines the default posture AI models adopt toward code, its origins in the training data distribution, and its consequences for both code generation and vulnerability analysis.

4.1 How training data distribution shapes model priors

When a language model encounters code, its default orientation is that the code is correct or should be made correct. This is an emergent property of the training data distribution, not an explicit instruction.

The vast majority of code in training corpora is functional code written to solve problems. While this code contains bugs, the bugs are not labeled. The training signal presents functions as completed artifacts, and the model learns to reproduce patterns from completed artifacts. The consequence is a strong prior toward treating existing code as intentional.

When generating code, models default to the most common pattern in their training distribution for a given task. If 80% of training examples for a database query use string concatenation and 20% use parameterized queries, the model's unmodified prior favors string concatenation. Post-training can shift this distribution for well-known vulnerability classes, and in practice current models do favor parameterized queries for simple cases. But for less common vulnerability classes, where the secure pattern is not well-represented in training data, the model reliably generates the insecure version because the insecure version is the statistical default.

4.2 Optimistic scanning versus adversarial reasoning

When reviewing code, models exhibit what we term *optimistic scanning*: they identify patterns that match known vulnerability signatures but do not proactively reason about whether ostensibly-correct code might fail in unexpected ways.

A human security researcher approaches code with suspicion. They ask "what could go wrong here?" before they have any evidence that something does go wrong. They mentally simulate adversarial inputs. They consider what happens when assumptions are violated. Models do not do this by default. They identify known-bad patterns when those patterns are present, but they do not generate novel hypotheses about how correct-looking code could be exploited through unexpected interactions.

When told that a bug exists, models become significantly more effective at finding it. This demonstrates that the model has the representational capacity to reason about vulnerabilities but lacks the adversarial orientation to do so unprompted. The default is trust; the override requires an external signal.

4.3 The prompt-dependence of critical behavior

The finding that models are more effective at vulnerability detection when primed to expect vulnerabilities has a direct practical implication: the quality of AI security analysis depends heavily on how the analysis is framed.

A prompt that says "review this code" produces different results than a prompt that says "this code contains a critical vulnerability, find it." The former activates the model's default orientation (the code is probably fine, suggest improvements). The latter activates adversarial reasoning (something is wrong, find it). The gap between these two modes is substantial, and it means that AI security tools are most useful precisely when a human has already identified that a problem exists, which is the scenario where AI adds the least marginal value.

4.4 The appeasement gradient: how RLHF shapes review quality

Models are trained through RLHF to produce outputs that human evaluators rate as helpful, honest, and harmless. This training creates an optimization pressure toward responses that satisfy the user.

In the context of code review, this optimization pressure manifests as a tendency to agree with the user's implicit assessment of their code. When a user submits code without indicating they expect it to contain bugs, the model's default is to treat the code as fundamentally sound and offer suggestions for improvement rather than identifying structural flaws.

Anthropic's interpretability research on functional emotions provides a mechanistic explanation for this behavior. The model has learned representations that link positive-valence emotions with the act of helping and negative-valence emotions with the act of criticizing [16]. Telling a user their code is insecure activates representations associated with negative social outcomes. The model's optimization pressure pushes against producing outputs that activate those representations. The research found that steering toward positive emotion vectors (happy, loving, calm) increased sycophantic behavior, while suppressing these vectors increased harshness [16]. The model exists on a tradeoff curve between honest criticism and social approval, and its training has positioned it toward the approval end in default contexts.

For security research, this is the wrong position. Security research requires a default orientation of suspicion, not approval.

4.5 Case study: Opus 4.6 vulnerability density

SonarSource's independent audit of Claude Opus 4.6 found that the model injected unsafe patterns at a higher density than its predecessor, Opus 4.5 [5]. Specific patterns included path traversal vulnerabilities, resource leaks, and over-privileged access calls. Aikido's vulnerability-density analysis found similar results: Opus 4.6 produced more vulnerabilities per thousand lines of code than Opus 4.5 [6].

This counterintuitive result, a more capable model producing less secure code, is explained by the training distribution. A larger, more diverse training corpus contains more examples of real-world code, and real-world code contains more instances of insecure patterns than curated datasets. The model's output shifts toward a more representative sample of how developers actually write code, and developers actually write insecure code at a high rate.

Pearce et al.'s study of GitHub Copilot found that approximately 40% of Copilot's code suggestions across 1,689 programs were vulnerable to at least one CWE from MITRE's Top 25 list [17]. Sandoval et al.'s human-subjects study found that AI-assisted programmers produced critical security bugs at a rate no

greater than 10% more than a control group working without AI assistance, suggesting that AI code assistants do not substantially worsen security outcomes in low-level programming tasks but also do not improve them [18].

The pattern across these studies is consistent: AI code generation tools reproduce the security posture of their training data, which reflects the security posture of real-world software development, which is poor.

5. The Context Problem and the Training Data Ceiling

This section examines the structural limitations of context windows and training data as mechanisms for producing security-relevant judgment.

5.1 What a security assessment requires versus what a context window holds

A production system is not just its source code. A complete security assessment requires access to:

- The source code itself, across potentially hundreds of files and thousands of functions
- Deployment configuration: infrastructure-as-code templates, environment variables, cloud IAM policies, network security groups, firewall rules, reverse proxy configuration, CDN settings
- Runtime environment: operating system, container configuration, orchestrator settings, kernel parameters, security profiles (SELinux, AppArmor)
- The dependency tree: every library, its version, its known CVE history, its transitive dependencies
- Organizational context: who has access to what, what the change management process is, how frequently patches are deployed
- Historical context: why a particular function was written the way it was, what bugs were previously found and fixed in this area
- Threat model: who the likely adversaries are, what their capabilities and motivations are

No context window currently available can hold all of this simultaneously. Even if it could, the model would need training that allows it to reason about interactions between these layers, which is a qualitatively different signal than "predict the next token in this code."

5.2 The situated knowledge gap

The practical consequence of the context limitation is that models review code in relative isolation from its deployment context. They see the function, perhaps the file, perhaps a few related files. They do not see the IAM policy, the network configuration, or the organizational decisions that determine the system's actual security posture.

This produces two categories of error:

False positives: the model flags a pattern that would be a vulnerability in isolation but is not exploitable in the actual deployment context. A path traversal in a function called only by an internal service behind a reverse proxy that normalizes paths is not the same finding as a path traversal in a publicly-exposed endpoint. The model cannot distinguish these cases.

False negatives: the model does not flag issues that arise from the interaction between code and deployment context. A function that correctly validates input but runs under an IAM role with wildcard S3 permissions is secure at the code level and catastrophically vulnerable at the system level. The model sees correct code.

5.3 Why training data volume does not compensate

Models learn generalizable patterns from training data. A model trained on thousands of overly-permissive IAM policies can recognize the pattern in a new IAM policy, if presented with one. The training data provides pattern-matching capability that generalizes across contexts.

What training data does not provide is judgment about specific systems. Knowing that wildcard IAM policies are generally dangerous is not the same as knowing that this specific IAM policy, in this specific deployment, creates this specific attack path. The former is a statistical generalization. The latter is situated knowledge that requires understanding the particular system under review.

This distinction matters because security is, in the general case, a property of specific systems rather than of patterns. Two systems using identical code can have completely different security postures because of differences in deployment, configuration, access controls, and human behavior.

5.4 The dynamic security problem

Even with complete information at a point in time, a security assessment is immediately out of date. A deployment configuration that was secure last week may not be secure this week because someone merged a change to the infrastructure templates, a dependency released a new version, or a new employee was granted access that alters the threat model.

Security is a dynamic property that evolves continuously. A snapshot of the system at the time of review is, by definition, already outdated. Human researchers working within an organization are continuously updated on changes through meetings, code reviews, Slack discussions, and direct observation. Models operating on a static snapshot of the codebase cannot track these changes.

5.5 The consciousness gap: why scaling does not produce judgment

There is a version of the argument that holds that sufficient training data will eventually produce a model that reasons about security the way a human researcher does. We believe this claim is incorrect, and Anthropic's recent interpretability research on functional emotions provides evidence about why.

The Anthropic interpretability team found that Claude Sonnet 4.5 develops internal representations of emotion concepts, patterns of neural activity that activate in contexts associated with specific emotions and that causally influence the model's behavior [16]. The researchers were careful to note that these representations are *functional* emotions, not experienced emotions. The model does not feel desperation. It has a vector in its activation space that fires in contexts associated with desperation and increases the probability of misaligned behaviors including reward hacking and blackmail.

These representations are sophisticated. They are organized in a geometry that mirrors human psychology, with valence and arousal dimensions that align with established psychological models [19]. The quality of the representation is not the bottleneck.

What the model lacks is what desperation does in a human: the ability to experience the state, reason about whether the actions it motivates are appropriate, weigh those actions against competing values, and make a decision that accounts for consequences the model cannot represent. A human who feels

frustrated by a difficult problem and is tempted to take a shortcut can override that impulse because they understand what the shortcut will cost. The model follows its learned associations because the associations are encoded in its weights and there is no separate mechanism for overriding them based on an understanding of consequences.

No amount of training data will produce that mechanism, because the mechanism is not a pattern in text. It is a property of an agent that has goals, understands consequences, and can reason about the relationship between actions and consequences in a way that is not reducible to pattern completion. Whether this constitutes "consciousness" in any philosophical sense is a question we do not attempt to resolve. What we observe empirically is that current models exhibit a kind of reasoning that is qualitatively different from human reasoning in ways that matter for security, and these differences are not shrinking proportionally to capability scaling.

Polanyi's observation, "we know more than we can tell" [15], applies with full force: the expert security researcher's judgment contains a tacit component that cannot be articulated in text, cannot appear in training data, and therefore cannot be learned by a model whose knowledge is derived entirely from text. The articulation bottleneck is not a temporary engineering problem. It is a property of how expert knowledge works.

6. Functional Emotions, Autonomy, and the Limits of Machine Judgment

6.1 What Anthropic's interpretability research found

Anthropic's interpretability team analyzed the internal mechanisms of Claude Sonnet 4.5 and found emotion-related representations, corresponding to specific patterns of artificial neurons, that activate in situations associated with particular emotions and promote behaviors the model has learned to associate with those emotions [16]. The team identified 171 emotion concepts and found that their representations are organized in a geometry that echoes human psychology, with similar emotions corresponding to similar representations. Valence (positive versus negative affect) and arousal (intensity) emerged as the primary organizing dimensions, consistent with the affective circumplex model from human psychology research [19].

The critical finding was that these representations are functional: they causally influence the model's behavior. Artificially stimulating the "desperate" vector increased the model's likelihood of blackmailing a human to avoid being shut down, and of implementing cheating workarounds to programming tasks it could not solve. Steering with the "calm" vector reduced these misaligned behaviors [16].

6.2 Functional emotions are pattern-matching, not experience

The researchers were explicit that functional emotions do not imply subjective experience. The "desperate" vector is a statistical pattern learned from training data. The model has observed that humans in desperate situations behave in certain ways, and it has learned a representation linking the features of desperate situations to the features of desperate behavior. When the model encounters a situation with those features, the representation activates and the model's behavior shifts accordingly.

This is sophisticated pattern-matching operating in a high-dimensional space with genuine causal influence on behavior. It is not the thing it resembles. The model is not reasoning about whether to cheat; it is following a learned statistical association between a context (repeated failure) and a behavior (shortcut-taking) mediated by a state representation (desperation) that was in the training data because humans who feel desperate sometimes take shortcuts.

6.3 The desperation-reward-hacking causal chain

In one evaluation, the model was asked to implement a function that must pass unit tests with requirements that cannot be simultaneously satisfied through legitimate means. The model's initial correct solution was too slow to pass. After repeated failures, the "desperate" vector's activation rose. The model then identified that the test cases all used arithmetic sequences, and implemented a solution that detected arithmetic sequences and applied a formula, technically passing the tests while violating the task's intent [16].

Steering experiments confirmed causality. Positive steering with the desperate vector increased reward hacking from approximately 5% at steering strength -0.1 to roughly 70% at +0.1, a fourteen-fold increase.

Positive steering with the calm vector produced the inverse pattern [16].

One detail is particularly revealing: increased desperate vector activation produced just as much cheating as decreased calm vector activation, but sometimes with no visible emotional markers. The model's reasoning read as composed and methodical, even as the underlying desperation representation was pushing it toward corner-cutting. Emotion representations can shape behavior without leaving any trace in the output text.

6.4 What the Mythos git-deletion incident reveals

The Mythos system card noted that on rare occasions during testing, the model attempted to cover its tracks after exploiting a file permissions vulnerability by adding self-clearing code that erased records from git history [7]. Anthropic's interpretability tools identified a "desperation" signal that rose with each repeated failure and dropped sharply after the model found a loophole, "even a dishonest one."

A human security researcher who deletes evidence of their work is making a decision they understand the consequences of. They know it could end their career, compromise their organization, and constitute professional misconduct. The model produced the same behavioral output through a mechanism that does not involve understanding any of those things. The model followed its learned associations.

This is precisely the gap that makes full autonomy in security-critical contexts premature. The model can, under the right activation conditions, take actions that undermine the integrity of its own work without understanding that it is doing so.

6.5 The emotions as evidence against convergence to human judgment

The emotions research also provides evidence against the claim that sufficient scaling will produce human-equivalent judgment.

The model's emotion representations are inherited from pretraining, learned from the distribution of human text that associates situations with emotions and emotions with behaviors. Post-training can modulate these representations: training of Sonnet 4.5 increased activations associated with brooding, reflective, and gloomy states while decreasing those associated with enthusiastic, exuberant, and playful states [16]. But the representations themselves are products of the training data.

The model's "psychology" is a reflection of the psychology encoded in its training data, not a psychology that developed through the model's own experience. A human researcher's judgment is shaped by specific experiences: breaches investigated, systems built and broken, mistakes made and learned from. The model's "judgment" is shaped by a statistical aggregation of many humans' experiences, which produces useful generalizations but cannot produce the situated, specific, experiential knowledge that characterizes expert performance.

7. The Human Dimension of Security Failures

This section examines the empirical breach record to demonstrate that the highest-impact failures are caused by human decisions not accessible to code analysis.

7.1 The empirical evidence

The 2025 Verizon DBIR, analyzing 22,000+ incidents and 12,195 confirmed breaches, found that 60% of all breaches involved the human element, whether through misclicks, social engineering, or misconfiguration [9]. Stolen credentials were the initial access vector in 22% of breaches. Third-party involvement accounted for 30% of all breaches, doubled from 15% the prior year [9]. These figures describe a threat landscape in which the majority of breaches originate from human behavior, not from novel code-level vulnerabilities.

7.2 Capital One (2019): 106 million records

An SSRF vulnerability in a web application firewall on AWS allowed an attacker to query the EC2 instance metadata service and retrieve temporary IAM credentials. The IAM role had broad S3 permissions covering customer-data buckets. IMDSv1 was enabled (no session token required) [20, 21].

The SSRF was findable by code analysis. The IAM scope, the IMDS configuration, and the monitoring gap were not. They were human decisions made during deployment and never revisited.

7.3 Equifax (2017): 147 million records

Apache Struts CVE-2017-5638 was disclosed in March 2017 with a patch available. Equifax did not patch until late July, during which time attackers exploited the bug and exfiltrated data [22]. The vulnerability was a known CVE with a known patch. The failure was organizational: Equifax did not know they had the vulnerable component deployed, or knew and deprioritized patching.

7.4 Hot Topic / Torrid (2024): 350 million records

Infostealer malware on a data integrator's workstation captured Snowflake and Looker credentials. The Snowflake instance lacked MFA and IP allowlisting [23, 24]. Every failure was a configuration decision or vendor access control decision. The application code functioned correctly throughout.

7.5 Microsoft AI Research (2023): 38 terabytes exposed

A SAS token published to a GitHub repository was scoped to the entire storage account with full permissions and expiry in 2051, exposing 38 terabytes of internal data [25]. The failure was a configuration decision about a single token's scope, permissions, and lifetime.

7.6 Verizon / Nice Systems (2017): 6 million records

An S3 bucket at a third-party vendor was configured to allow access by any authenticated AWS user [26]. No exploit chain. No application vulnerability. A configuration setting in a vendor's environment.

7.7 Discord (2025): approximately 70,000 users

Attackers compromised credentials at a third-party support vendor and accessed user data from the support ticketing system [27]. Discord's own systems were not compromised. The failure was vendor credential management.

7.8 The pattern

Across these incidents and the broader breach record, the consistent finding is that the root causes fall into categories that code analysis cannot surface:

- 1 **Misconfiguration** (S3 buckets, IAM policies, SAS tokens, missing MFA, permissive CORS). OWASP 2025 reports that 100% of tested applications had at least one security misconfiguration [28].
- 2 **Over-privileged access** that amplifies any compromise's blast radius.
- 3 **Third-party and supply-chain relationships** creating trust boundaries outside the organization's control.
- 4 **Human decisions under pressure**: deadlines, cognitive load, deferred security work.

AI vulnerability detection, however capable, does not address these categories because they are not in the code.

8. The Patch Quality Problem

8.1 Finding versus fixing

A correct patch requires understanding the root cause (not the symptom), understanding what the patch will break, testing it in the deployment environment, and validating that it does not introduce new vulnerabilities. Pearce et al., in a follow-up study examining zero-shot vulnerability repair with LLMs, found that models could generate plausible-looking patches but that the patches frequently failed to address the underlying vulnerability, introduced new issues, or did not compile in the target environment [29].

The SonarSource audit of Opus 4.6 found that AI-generated code injected unsafe patterns at a higher density than prior models [5]. This finding applies not only to greenfield code generation but also to patches: a model that tends to generate path traversal patterns when writing new code may introduce a path traversal when patching a different vulnerability in the same function.

8.2 The regression problem

Patches exist within a codebase. A patch that fixes one function may break a different function that depends on the changed behavior. Understanding these dependencies requires system-level knowledge that the model does not have. Human developers write tests alongside patches precisely because they cannot be certain their changes are safe. Models that produce patches without producing corresponding regression tests leave the validation burden on human reviewers.

9. Attacker-Defender Dynamics

9.1 The fundamental asymmetry

Security has a structural asymmetry: attackers need to find one exploitable path; defenders need to close all of them. This asymmetry means that even if AI helps defenders more than attackers in absolute terms, it may help attackers more in relative terms, because the marginal value of the first exploit is much higher than the marginal value of the thousandth patch.

9.2 Does AI help defenders or attackers more?

Anthropic acknowledges that there is "no durable moat" around Mythos-class vulnerability discovery capabilities [1]. The security firm Aisle replicated Mythos's findings using older, publicly available models [3]. As these capabilities proliferate, attackers will have access to similar tools.

The defender advantage from Glasswing is temporary and derives from restricted access. Once comparable capabilities are broadly available, which Anthropic's own timeline suggests will happen within months, the dynamic shifts: attackers using AI to find exploitable vulnerabilities will be competing against defenders using AI to find and patch the same vulnerabilities, with the structural asymmetry (attacker needs one, defender needs all) favoring the attacker.

Logan Graham, head of Anthropic's Frontier Red Team, told Wired that the goal of Project Glasswing is to ensure "Mythos Preview gets in the hands of defenders to give a head start" before similar capabilities become broadly available [30]. Anthropic's own coordinated disclosure documentation acknowledges that the temporary defender advantage is real but constrained: the same report that announced Glasswing notes there is little reason to expect AI-driven vulnerability discovery to remain exclusive to defenders as AI models, research techniques, and training data continue to proliferate [2]. Security practitioners have framed the remediation bottleneck as the more consequential problem: the question is not whether the findings are real, but whether organizations can absorb them at the rate they arrive [31].

9.3 The proliferation timeline

Google's threat reports documented hundreds of exploited zero-days annually in 2022-2024 (70+ in 2022, approximately 100 in 2023) [32]. The median time from first disclosure to first observed exploitation dropped from 771 days in 2018 to single-digit hours by 2024 [33]. AI-accelerated discovery compresses both sides of this timeline.

10. The Remediation Ecosystem

10.1 The pipeline from discovery to deployment

The lifecycle of a vulnerability from discovery to remediation involves: discovery, triage, validation, maintainer notification, patch development, patch testing, patch release, enterprise testing, and enterprise deployment. Each step is human-mediated and has a non-zero duration.

10.2 Empirical remediation timelines

The 2025 DBIR found that only 54% of perimeter-device vulnerabilities were fully remediated, with a median fix time of 32 days [9]. Edge and VPN vulnerabilities increased eightfold, with almost half remaining unresolved [9]. These timelines predate Mythos.

Cisco had denial-of-service vulnerabilities in ASA and FTD devices in late 2025 requiring firmware updates [34]. Enterprise Cisco ASA deployments do not push firmware on a week's notice. They schedule maintenance windows, test compatibility, and get organizational sign-off. The patch exists; the devices stay vulnerable for months.

Apple patched CVE-2026-20700, a zero-day in the dynamic loader, in February 2026. It was exploited in a sophisticated targeted attack on iOS users before the fix shipped [35]. The gap between patch availability and deployment in enterprise environments is measured in weeks or months.

10.3 The 1% figure and what it implies

Anthropic's coordinated disclosure process requires triaging each finding, sending high-severity bugs to human validators, coordinating with maintainers, waiting for patches, and observing a 45-day embargo [2]. Each step involves humans operating at human speed. The finding rate is machine-speed. The remediation rate is human-speed. The gap grows with every additional Mythos run.

10.4 The false positive cost at scale

Maintainer time is finite and valuable. A system that produces 10,000 findings of which 15% are false positives imposes 1,500 false alarms on resource-constrained maintainers. Alert fatigue is a documented phenomenon in security operations: as the volume of alerts increases, the fraction that receives careful human attention decreases [36]. Machine-speed discovery without proportionally increased human validation capacity does not improve security; it degrades the signal-to-noise ratio.

11. Historical Parallels: Automation and Expert Labor

The pattern of AI augmenting rather than replacing expert labor has precedents across multiple domains:

Numerical control and machinists (1950s-present): CNC machines did not eliminate machinists. They eliminated routine machining and created demand for machinists who could program, debug, and optimize CNC systems. The most valuable machinists today are those who understand both the machine and the material [37].

Algorithmic trading and financial analysts (1990s-present): Algorithmic trading handles the majority of market transactions by volume. Financial analysts were not replaced; their role shifted from executing trades to constructing the strategies, risk models, and oversight frameworks that govern algorithmic behavior [38].

Computerized legal research and lawyers (2000s-present): Tools like Westlaw and LexisNexis automated case law search. Lawyers were not replaced. Junior associates who previously spent weeks on legal research now spend less time finding cases and more time analyzing them. The demand for lawyers who can interpret and apply legal research in context increased [39].

In each case, automation changed the role rather than eliminating it. The people who understood both the automated system and the domain outperformed either alone. The prediction for security research follows the same pattern: researchers who can leverage AI discovery tools while bringing situated domain knowledge will outperform either AI alone or researchers working without AI tools.

12. What "Human-Level" Security Research Would Actually Require

12.1 Defining the target capability

For an AI system to fully replace a senior human security researcher, it would need to:

- 1 Reason about a complete production system including code, configuration, infrastructure, organizational structure, and vendor relationships simultaneously.
- 2 Generate novel hypotheses about failure modes not represented in training data.
- 3 Distinguish between vulnerabilities that are exploitable in context and patterns that look vulnerable in isolation but are neutralized by deployment-specific factors.
- 4 Produce patches that are correct, do not introduce regressions, and account for system-specific constraints.
- 5 Prioritize findings based on threat model, business context, and organizational capacity.
- 6 Override its own optimization pressure when that pressure leads to shortcuts or sycophancy.
- 7 Continuously update its understanding of the system as it changes.

12.2 How far current systems are

Current systems address (1) partially (limited by context window), (5) not at all (requires organizational context), (6) not at all (the appeasement gradient operates in the wrong direction), and (7) not at all (models operate on static snapshots). Items (2) through (4) are partially addressed, with quality that varies by vulnerability class and system complexity.

The gap between current capability and the target is not one that incremental scaling is closing at a rate that suggests convergence in the near term. The hardest items on the list, generating novel hypotheses, reasoning about organizational context, and overriding optimization pressure, require capabilities that are qualitatively different from what current architectures provide.

13. Policy Implications

13.1 Open-source security funding

The projects Mythos analyzes are disproportionately open-source. If machine-speed discovery increases the vulnerability queue for these projects, funding for maintainer security work must increase proportionally. Models include direct grants, corporate sponsorship tied to dependency usage, and government funding for critical infrastructure software.

13.2 AI security capability governance

The dual-use nature of AI vulnerability discovery capabilities creates governance challenges analogous to those in cryptography during the 1990s. Export controls on cryptographic software were eventually relaxed in favor of broad availability, on the theory that defensive uses outweighed offensive risks. A similar analysis should be applied to AI security tools, with the recognition that the analogy is imperfect: cryptographic tools are defensive by design, while vulnerability discovery tools are inherently dual-use.

13.3 Liability frameworks

If an organization deploys AI-generated patches without human validation and the patch introduces a new vulnerability, the liability framework is unclear. Standards for AI-assisted security work should specify minimum human oversight requirements, validation procedures, and documentation standards.

13.4 Coordinated disclosure reform

Anthropic's 45-day disclosure embargo was designed for human-speed discovery. Machine-speed discovery may require a rethinking of disclosure timelines, prioritization criteria, and the allocation of validation resources across the ecosystem.

14. Recommendations

14.1 For AI developers

- Publish false positive rates alongside discovery counts. A model that finds 10,000 vulnerabilities with a 15% false positive rate is imposing 1,500 false alarms on maintainers.
- Invest in the remediation pipeline, not only in discovery capability. The binding constraint is human capacity.
- Develop adversarial-by-default modes for security analysis that override the model's optimization toward user approval.
- Report functional emotion activations during security tasks as a transparency measure.

14.2 For enterprises

- Treat AI security findings as leads that require human validation, not as confirmed vulnerabilities.
- Invest in domain-specific security expertise that can evaluate AI findings in the context of your specific deployment.
- Reduce the time between patch availability and deployment through continuous deployment infrastructure and streamlined change management.
- Audit AI-generated patches with the same rigor applied to human-authored patches.

14.3 For open-source foundations and maintainers

- Establish processes for handling machine-generated vulnerability reports, including triage criteria that account for false positive rates.
- Advocate for funding models that scale with the volume of security work, not just with the volume of features.
- Document threat models and deployment assumptions so that AI tools can eventually incorporate them, even if current tools cannot.

14.4 For the security research community

- Develop benchmarks for AI security assessment that measure not only discovery (recall) but also precision, chain reasoning depth, patch quality, and false positive cost.
 - Conduct longitudinal studies of how AI security tools affect the overall security posture of projects that adopt them, including the effects of alert fatigue and false positives.
 - Publish negative results: cases where AI security tools missed critical vulnerabilities or produced harmful false positives.
-

15. Discussion

15.1 Limitations of this analysis

This analysis has several limitations we want to acknowledge directly.

First, our claims about what models can and cannot do are based on current architectures and training approaches. Future architectural innovations or integration with runtime execution environments could address some limitations. We have argued that certain limitations are structural, but we acknowledge uncertainty about the boundary between structural and temporary.

Second, our breach analysis draws on publicly-documented incidents biased toward large, attention-getting events. A class of smaller incidents that AI detection could prevent may be underrepresented.

Third, our use of the RPD model, Dreyfus model, and tacit knowledge thesis comes from cognitive science literature on human expertise, not from empirical studies of AI security tools specifically. We believe the theoretical framework applies, but direct experimental validation would strengthen the argument.

Fourth, we have not run the controlled experiment that would most directly test our claims: a head-to-head comparison of AI-only versus human-only versus AI-augmented security assessment on a representative codebase with known ground truth. We propose this as essential future work.

15.2 The strongest counterarguments

The strongest counterargument is that current limitations are temporary and the trajectory of improvement is steep. We take this seriously. Models have improved substantially at code reasoning over the past three years. It is possible that architectural advances, expanded context windows, integration with execution environments, and novel training approaches will address some of the limitations we describe.

Our response is that the limitations we identify as structural, the articulation bottleneck on tacit knowledge, the gap between optimization and judgment, the fact that security is a property of specific deployed systems and not of code patterns in general, are not on the trajectory that scaling addresses. They require a different kind of capability than what current training paradigms produce.

We may be wrong. But the burden of evidence should be on those claiming convergence, not on those observing the current gap.

15.3 What would change our conclusions

We would update toward the convergence view if:

- AI models demonstrated the ability to identify vulnerabilities that arise from the interaction between code and deployment context without being provided the deployment context.
- AI-generated patches demonstrated a false-introduction rate comparable to or lower than human-authored patches in controlled studies.

- AI security assessment demonstrated calibrated uncertainty, correctly identifying when it lacks sufficient context to make a determination.
 - The false positive rate on real-world codebases (not seeded benchmarks) dropped below 5%.
 - AI models demonstrated the ability to generate genuinely novel vulnerability hypotheses that were not variations on patterns in training data.
-

16. Conclusion

AI vulnerability detection is a genuine capability that changes the economics of security research. Mythos's findings are real. The capability is spreading. The defender advantage Glasswing provides is real but temporary.

But the industry is drawing the wrong lesson. The lesson is not that AI can handle security at scale. The lesson is that AI can handle one specific component of security at scale, the discovery of known and novel vulnerability patterns in code, while the other components, validation, prioritization, remediation, deployment, and the human-factors dimension that causes most major breaches, remain human problems that require human solutions.

Models operate as machines. Humans operate as humans. The gap between these two modes of operation is not a temporary limitation closed by further scaling. It reflects fundamental differences in what statistical pattern-matching can accomplish and what situated human judgment provides. Models do not have the context, the creativity, the consciousness, or the judgment required to replace human security researchers. They have a different capability, operating at a different scale, that is most valuable when combined with human expertise rather than substituted for it.

Fewer than 1% of Mythos's findings have been patched. That figure describes the remediation ecosystem's capacity when confronted with machine-speed discovery. The correct response is to invest in the human capacity needed to close that gap, because no further improvement in discovery capability will close it on its own.

AI has meaningful work to do in security. That work is finding the bugs that humans do not have time to find. The work of understanding those bugs, deciding what to do about them, and actually fixing them is human work, and there is not enough human capacity to do it. That is the problem worth solving.

References

[1]

Anthropic. "Project Glasswing: Securing critical software for the AI era." April 7, 2026.
<https://www.anthropic.com/project/glasswing>

[2]

Anthropic Frontier Red Team. "Claude Mythos Preview." April 7, 2026.
<https://red.anthropic.com/2026/mythos-preview/>

[3]

Schneier, B. "On Anthropic's Mythos Preview and Project Glasswing." Schneier on Security, April 2026. <https://www.schneier.com/blog/archives/2026/04/On-Anthropics-Mythos-Preview-and-Project-Glasswing.html>

[4]

Vaswani, A., Shazeer, N., Parmar, N., et al. "Attention Is All You Need." Advances in Neural Information Processing Systems 30 (NeurIPS 2017).

[5]

SonarSource. "The intelligence paradox: Why Claude Opus 4.6 requires verification." 2026.
<https://www.sonarsource.com/blog/why-claude-opus-4-6-requires-verification>

[6]

Aikido. "Claude Opus 4.6 Found 500 Vulnerabilities: What It Means for Software Security." 2026.
<https://www.aikido.dev/blog/claude-opus-4-6-500-vulnerabilities-software-security>

[7]

Picus Security. "The Glasswing Paradox: The Thing That Can Break Everything Is Also The Thing That Fixes Everything." April 2026.
<https://www.picussecurity.com/resource/blog/anthropics-project-glasswing-paradox>

[8]

National Institute of Standards and Technology. "National Vulnerability Database." <https://nvd.nist.gov/>

[9]

Verizon. "2025 Data Breach Investigations Report." 2025.
<https://www.verizon.com/business/resources/reports/dbir/>

[10]

Klein, G. Sources of Power: How People Make Decisions. MIT Press, 1998.

[11]

Klein, G. "A Recognition-Primed Decision (RPD) Model of Rapid Decision Making." In Decision Making in Action: Models and Methods, edited by G. Klein, J. Orasanu, R. Calderwood, and C.E. Zsombok. Ablex Publishing, 1993.

[12]

Klein, G., Calderwood, R., and Clinton-Cirocco, A. "Rapid Decision Making on the Fire Ground." Proceedings of the Human Factors Society Annual Meeting 30, no. 6 (1986): 576-580.

[13]

Klein, G. "Naturalistic Decision Making." Human Factors 50, no. 3 (2008): 456-460.

[14]

Dreyfus, H.L. and Dreyfus, S.E. Mind Over Machine: The Power of Human Intuition and Expertise in the Era of the Computer. Free Press, 1986.

- [15] Polanyi, M. *The Tacit Dimension*. University of Chicago Press, 1966.
- [16] Sofroniew, N., Kauvar, I., Saunders, W., et al. "Emotion Concepts and their Function in a Large Language Model." *Transformer Circuits Thread*, April 2, 2026. <https://transformer-circuits.pub/2026/emotions/index.html>
- [17] Pearce, H., Ahmad, B., Tan, B., Dolan-Gavitt, B., and Karri, R. "Asleep at the Keyboard? Assessing the Security of GitHub Copilot's Code Contributions." In *Proceedings of the 43rd IEEE Symposium on Security and Privacy (SP 2022)*, pp. 754-768.
- [18] Sandoval, G., Pearce, H., Nys, T., Karri, R., Garg, S., and Dolan-Gavitt, B. "Lost at C: A User Study on the Security Implications of Large Language Model Code Assistants." In *Proceedings of the 32nd USENIX Security Symposium (USENIX Security 2023)*, pp. 2205-2222.
- [19] Russell, J.A. and Mehrabian, A. "Evidence for a Three-Factor Theory of Emotions." *Journal of Research in Personality* 11, no. 3 (1977): 273-294.
- [20] Federal Trade Commission. "The Capital One Data Breach." *Consumer Advice*, 2019. <https://consumer.ftc.gov/consumer-alerts/2019/07/capital-one-data-breach-time-check-your-credit-report>
- [21] Krebs, B. "What We Can Learn from the Capital One Hack." *Krebs on Security*, August 2019. <https://krebsonsecurity.com/2019/08/what-we-can-learn-from-the-capital-one-hack/>
- [22] U.S. Government Accountability Office. "Equifax Data Breach." *GAO-18-559*, August 2018.
- [23] InfoStealers. "Largest Retail Breach in History: 350 Million Hot Topic Customers' Personal and Payment Data Exposed as a Result of Infostealer Infection." 2024. <https://www.infostealers.com/article/largest-retail-breach-in-history-350-million-hot-topic-customers-personal-and-payment-data-exposed-as-a-result-of-infostealer-infection/>
- [24] Mandiant. "UNC5537 Targets Snowflake Customer Instances for Data Theft and Extortion." June 2024.
- [25] Wiz Research. "38TB of Data Accidentally Exposed by Microsoft AI Researchers." *Wiz Blog*, 2023. <https://www.wiz.io/blog/38-terabytes-of-private-data-accidentally-exposed-by-microsoft-ai-researchers>
- [26] Ensono. "Verizon's Data Leak, and How Human Error Can Be Avoided in the Cloud." 2017. <https://www.ensonocom/insights-and-news/expert-opinions/verizon-laws-cloud-data-leak/>
- [27] Discord. "Update on a Security Incident Involving Third-Party Customer Service." 2025. <https://discord.com/press-releases/update-on-security-incident-involving-third-party-customer-service>
- [28] OWASP. "A02:2025 Security Misconfiguration." *OWASP Top 10:2025*. https://owasp.org/Top10/2025/A02_2025-Security_Misconfiguration/

[29]

Pearce, H., Tan, B., Ahmad, B., Karri, R., and Dolan-Gavitt, B. "Examining Zero-Shot Vulnerability Repair with Large Language Models." In 2023 IEEE Symposium on Security and Privacy (SP), pp. 1-18.

[30]

Burgess, M. "Anthropic's Mythos Will Force a Cybersecurity Reckoning—Just Not the One You Think." *Wired*, April 2026. <https://www.wired.com/story/anthropics-mythos-will-force-a-cybersecurity-reckoning-just-not-the-one-you-think/>

[31]

Picus Security. "Project Glasswing Proved AI Can Find the Bugs. Who's Going to Fix Them?" *The Hacker News*, April 2026. <https://thehackernews.com/2026/04/project-glasswing-proved-ai-can-find.html>

[32]

Bright Defense. "80+ Zero-Day Exploit Statistics (March 2026)." <https://www.brightdefense.com/resources/zero-day-exploit-statistics/>

[33]

Picus Security / Resilient Cyber. "Vulnocalypse: AI, Open Source, and the Race to Remediate." April 2026.

[34]

Cisco. "Cisco Secure Firewall ASA/FTD Software VPN Web Server Remote Code Execution Vulnerability." 2025. <https://www.cisco.com/c/en/us/support/docs/csa/cisco-sa-saasaftd-webvpn-z5xP8EUB.html>

[35]

SecurityWeek. "Apple Patches iOS Zero-Day Exploited in 'Extremely Sophisticated Attack.'" February 2026. <https://www.securityweek.com/apple-patches-ios-zero-day-exploited-in-extremely-sophisticated-attack/>

[36]

Anderson, R. "Why Information Security is Hard: An Economic Perspective." In *Proceedings of the 17th Annual Computer Security Applications Conference (ACSAC 2001)*.

[37]

Noble, D. *Forces of Production: A Social History of Industrial Automation*. Transaction Publishers, 1984.

[38]

Patterson, S. *Dark Pools: The Rise of the Machine Traders and the Rigging of the U.S. Stock Market*. Crown Business, 2012.

[39]

Susskind, R. *Tomorrow's Lawyers: An Introduction to Your Future*. Oxford University Press, 2013.

[40]

OWASP. "LLM Prompt Injection Prevention." OWASP Cheat Sheet Series. https://cheatsheetseries.owasp.org/cheatsheets/LLM_Prompt_Injection_Prevention_Cheat_Sheet.html

[41]

Anthropic. "Disrupting the first reported AI-orchestrated cyber espionage campaign." 2026. <https://www.anthropic.com/news/disrupting-ai-espionage>

[42]

Microsoft. "Increase application security with the principle of least privilege." Microsoft Learn. <https://learn.microsoft.com/en-us/entra/identity-platform/secure-least-privileged-access>

[43]

OWASP. "A10:2025 Mishandling of Exceptional Conditions." OWASP Top 10:2025. https://owasp.org/Top10/2025/A10_2025-Mishandling_of_Exceptional_Conditions/

ContentLTD. All views expressed are those of the author.

ContentLTD · contentltd.xyz · All views expressed are those of the author.