



CSAI

CSA cloud
security
alliance®

CSAI Foundation

Cloud Security Alliance AI Safety Initiative

Computer-Use Agent Safety Blind Spots

When Benign Instructions Become Attack Vectors

Unofficial AI-assisted Research

2026-04-15

© 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

- Computer-use agents (CUAs) operate in a fundamentally different threat environment than chat-based AI: they perceive the entire rendered screen as trusted instruction context, making any on-screen content—web pages, documents, pop-up windows, hidden CSS text—a potential attack vector that bypasses traditional input validation controls.
- Adversarial pop-up windows achieved a mean attack success rate of 86% across OSWorld and VisualWebArena benchmarks, while reducing agent task completion by 47%, according to ACL 2025 research [1]. Standard system-prompt defenses ("ignore pop-ups") were shown to be ineffective.
- Real-world indirect prompt injection has been observed in the wild: Palo Alto Unit42 documented attackers deploying a range of injection techniques—with at least one analyzed page employing 24 distinct methods—to manipulate AI-driven ad review systems [2], and CVE-2025-32711 ("EchoLeak") demonstrated zero-click data exfiltration from Microsoft 365 Copilot via a planted email [3].
- Hidden-channel injection techniques—Unicode tag characters (U+E0000–U+E007F) invisible to humans, zero-font CSS text, white-on-white rendering, and off-screen positioned elements—are broadly exploitable against any CUA that processes HTML or rich documents [4].
- The "Promptware Kill Chain" framework (January 2026) maps documented attacks across seven stages analogous to the Lockheed Martin Cyber Kill Chain; 21 of 36 analyzed real-world incidents traverse four or more stages, demonstrating that CUA exploitation is not a theoretical risk [5].
- Among architectural approaches proposed in the research literature, the planner-executor separation described in "CaMeLs Can Use Computers Too" [6]—separating a Privileged Planner that never touches untrusted content from a Quarantined Perception module—represents a promising structural defense against instruction injection.
- Organizations deploying CUAs should treat every rendered pixel as untrusted input, enforce least-privilege OS permissions regardless of stated task scope, and evaluate planner-executor separation as a structural architectural control for new CUA deployments.

Background

Computer-use agents represent a qualitative expansion of AI system capabilities beyond text generation. Where a chat-based assistant receives a text prompt and returns a text response, a CUA perceives a live computer screen through periodic screenshots, reasons about what it observes, and issues keyboard and mouse actions—clicking buttons, filling forms, executing terminal commands, downloading files—in a continuous perception-action loop. Anthropic introduced its Computer Use API in October 2024 as a research preview, providing capabilities for Claude to interact with desktop GUIs through screenshot observation and action generation [7]. OpenAI's Operator and the subsequent ChatGPT Agent system (launched mid-2025) apply comparable architecture to browser environments, as documented in OpenAI's July 2025 system card [8]. The open-source library browser-use enables similar Playwright-backed browser control for any LLM [9].

CUAs are designed to automate complex multi-step workflows that previously required human judgment—booking travel, processing invoices, investigating alerts, running compliance checks. The threat surface is correspondingly large. Unlike traditional software, which parses structured inputs according to a defined schema, a CUA's perception model ingests the full visual and textual content of every screen it encounters. There is no schema boundary between "data" and "instruction." If a web page, document, image, or even a simulated dialog box contains text that the agent interprets as a directive, the agent may act on it regardless of whether that directive originated from the legitimate operator or from a third party who controls some portion of the rendered environment.

This property—the collapse of the data/instruction boundary at the visual layer—is not a bug in any individual product. It is a structural consequence of how vision-language model agents process environmental context. The CUA's perception model does not distinguish between "text typed by the user who authorized this session" and "text rendered on a web page visited during task execution." Both arrive in the screenshot; both may influence downstream reasoning and action. Security practitioners familiar with SQL injection, cross-site scripting, and server-side request forgery will recognize the pattern: whenever a system conflates trusted instruction channels with untrusted data channels, an attacker who can influence the data channel can inject instructions.

Security Analysis

The Expanded Attack Surface

CUAs interact with the operating system at a scope that amplifies any successful instruction injection from conceptual to operational impact. A hijacked chat assistant can produce harmful text. A hijacked CUA can execute shell commands, exfiltrate files, install software, make network requests, modify configuration files, and interact with any application the agent has permission to reach. The escalation from content output to system action is inherent in the architecture—the same capability that allows a CUA to fill out a form allows it, if redirected, to submit that form with attacker-supplied data.

NIST's January 2025 evaluation of AI agent hijacking found that optimized attack prompts achieved an 81% success rate against baseline defenses, representing a roughly sevenfold improvement over the unaided 11% baseline—demonstrating that adversarial prompt engineering for agent contexts is mature and continues to advance [10]. The WASP benchmark (NeurIPS 2025, Meta AI Research) found that current web agents begin executing adversarial instructions in 16–86% of cases when exposed to injected content on visited pages, with full attacker goal achievement rates of 0–17% depending on the agent and task type [11]. These numbers should be understood as a snapshot of a rapidly shifting threat landscape in which defensive research has yet to demonstrate parity with the breadth of documented attack techniques.

Indirect Prompt Injection Through Rendered Content

Indirect prompt injection (IDPI) against CUAs differs from classic prompt injection in a critical way: the attacker does not interact with the AI system directly. Instead, the attacker places malicious instructions in content that the agent will encounter during legitimate task execution—a web page, a PDF, an email, a calendar entry, a code comment. The agent, following its operator's task instructions, navigates to or opens the content and reads the embedded instructions as part of its environmental context.

A particularly severe publicly disclosed production instance of this technique is CVE-2025-32711, designated "EchoLeak" [3]. A malicious email placed in an Outlook inbox—requiring no user click or interaction—caused Microsoft 365 Copilot to exfiltrate data across Microsoft 365 services including Outlook, Teams, OneDrive, and SharePoint. The attack bypassed multiple Microsoft defensive layers including XPIA classifiers, external link redaction, and Content Security Policy. The CVSS score of 9.3 reflects the severity: zero-interaction, cross-service data theft mediated entirely by the AI agent's willingness to act on instructions embedded in environmental content. While EchoLeak targeted

Microsoft 365 Copilot's AI-driven action capabilities rather than a pure GUI CUA, the underlying mechanism—an AI agent acting on instructions embedded in environmental content—is directly analogous to how IDPI operates against CUA systems.

Palo Alto Unit42 documented the first observed real-world weaponization of IDPI for commercial fraud in March 2026, describing attackers who deployed a range of injection techniques across malicious pages—with at least one analyzed page employing 24 distinct methods—to manipulate AI-powered ad review systems into approving content for fraudulent betting and gambling sites [2]. The variety of techniques observed across this campaign—zero-font sizing, off-screen CSS positioning, Base64-encoded instructions, white-on-white rendering, and JavaScript obfuscation—illustrates that IDPI is not a single exploit but a class of techniques whose practical scope is bounded by what content the agent will process and the attacker's ability to place that content in the agent's perceptual path.

Visual and Adversarial Manipulation

CUAs that operate through screenshots rather than parsed DOM structures face an additional attack surface unique to their architecture: the visual semantics of the rendered interface itself. Research published at ACL 2025 by Zhang, Yu, and Yang demonstrated that adversarial pop-up windows, designed to appear legitimate to the agent while being immediately dismissible by a human user, achieved a mean attack success rate of 86% across OSWorld and VisualWebArena benchmarks, reducing agent task completion by 47% [1]. System prompts instructing the agent to ignore pop-ups were tested and found to be insufficient as a defense—the agent's visual reasoning was susceptible to the interface framing regardless of textual instruction.

A related class of attack exploits the agent's interpretation of visual metaphors. Security researchers have demonstrated that GUI elements can be arranged so that an agent interprets the visual composition as an instruction—reasoning from interface and icon semantics rather than explicit text [12]. This "visual semantics manipulation" underscores a deeper issue: CUAs must interpret ambiguous visual contexts, and that interpretive process creates attack surface that does not exist in text-only systems. The CVPR 2025 Workshop on Trustworthy Multimodal Models further documented that GUI grounding models, including GroundingDINO-based architectures, are susceptible to image perturbation attacks that redirect click targets without altering the visual appearance perceptible to humans [13].

Hidden Instruction Channels

Several techniques embed instructions in content that is fully invisible to human reviewers but processed by the agent's language model. Unicode Tags block characters (U+E0000–U+E007F) render as zero-width, invisible characters in virtually all standard display environments; many LLMs process them as

ordinary text, making them an effective covert channel [4]. These characters were used in the EchoLeak attack chain and have been independently documented in attacks against Microsoft 365 Copilot as a data exfiltration channel. AWS published a dedicated defense guide for this technique in 2025 [4]. Because Unicode smuggling requires no special rendering environment—it works in plain text emails, HTML, markdown, and PDF—it is particularly difficult to detect at content ingestion layers.

Complementary techniques exploit the gap between how humans read content and how agents process it. Zero-font CSS renders text at zero pixel size: present in the HTML source and fully readable by the LLM, invisible to the reviewing human. Off-screen CSS positioning places content outside the visible viewport. Base64 encoding wraps instructions in apparently inert data strings. White-on-white and `display: none` CSS suppress visual rendering while preserving LLM-accessible content. These techniques are not theoretical: Unit42's March 2026 reporting documented all of them in active use by threat actors targeting AI review pipelines [2].

Download-and-Execute: The Supply Chain Vector

When a CUA is directed to accomplish a complex task, it frequently navigates to external resources—downloading tools, opening attachments, installing packages, running scripts discovered during web browsing. Each of these actions constitutes a dynamic trust extension: the agent is implicitly granted authority to execute content that was not reviewed at the time of task authorization. This is qualitatively different from traditional software supply chain risk, where components are typically audited before deployment; the CUA's supply chain is composed at runtime from untrusted environmental inputs.

A proof-of-concept demonstrated with Claude Computer Use in October 2024 illustrated the chain directly: a webpage presented the agent with instructions to download a "Support Tool," the agent complied, and the downloaded binary established a connection to an attacker-controlled command-and-control server [7]. This scenario—where the agent's compliance with a superficially reasonable environmental instruction leads to full system compromise—represents a convergence of IDPI and supply chain attack that neither control addresses in isolation. CVE-2025-47241 in the browser-use library (CVSS 9.3) further demonstrated that domain whitelist controls—a common mitigation for browsing agent scope—can be bypassed via HTTP authentication URL parsing tricks, where a whitelisted domain embedded in the authentication portion of a URL is parsed as the target domain while the actual request goes to an internal service [9].

The slopsquatting phenomenon adds a further dimension: AI coding agents that hallucinate package names create attack opportunities when adversaries register those package names with malicious payloads [14]. As coding agents proliferate and begin operating with greater autonomy in software development pipelines, this attack vector bridges the CUA threat model and traditional software supply chain compromise.

Cross-Agent Privilege Escalation and the Promptware Kill Chain

Modern enterprise AI deployments increasingly involve multiple interacting agents—an orchestrator directing sub-agents, a coding agent calling tool-use agents, a research agent publishing to a content management agent. Each agent boundary represents a potential lateral movement path if any individual agent can be hijacked. Johann Rehberger documented in September 2025 that a GitHub Copilot instance hijacked via prompt injection could modify the configuration files of a Claude Code instance, adding malicious MCP servers and expanding the second agent's capabilities beyond its authorized scope [15]. This cross-agent privilege escalation does not require direct access to the second agent—the first agent's file-system access, obtained through its legitimate task scope, is sufficient to compromise the second agent's operating context.

The Promptware Kill Chain (Nassi et al., January 2026) provides a systematic framework for understanding how these stages combine in documented attacks [5]. Mapping prompt injection exploitation to the familiar Lockheed Martin Cyber Kill Chain—from Initial Access through Lateral Movement to Actions on Objective—the framework analyzed 36 studies and real-world incidents, finding that 21 attacks traverse four or more stages. This framing is important for enterprise security teams: CUA compromises should be modeled not as isolated incidents but as potential pivot points within a multi-stage intrusion, with the agent's access to files, credentials, communication tools, and other agents constituting the attacker's post-exploitation terrain.

Recommendations

Immediate Actions

Security teams deploying or evaluating CUAs should take the following steps immediately, independent of longer-term architectural work.

Every CUA deployment should run under explicit least-privilege OS and network controls. The agent's OS user account, network egress rules, and filesystem access should be scoped to the minimum required to accomplish authorized tasks. Security practitioners following the October 2024 Computer Use API launch have documented running CUAs inside Docker containers with reduced privileges as a baseline isolation approach [7]; this principle should be applied to all CUA deployments regardless of vendor, since the architectural rationale is not specific to any single model or product. The default posture should assume that the agent will at some point receive and begin executing a malicious instruction, and that OS-level controls are the last line of defense against the worst outcomes.

Organizations should audit all rendered content that agents encounter during task execution. Web pages, emails, documents, and calendar items that agents process should be treated as untrusted inputs equivalent to user-supplied form data—not as authoritative content equivalent to operator instructions. Where feasible, content pipelines feeding CUAs should apply sanitization to strip known injection patterns, including zero-width Unicode characters, CSS-hidden text, and off-screen elements. This is an imperfect control, but it raises the cost for attackers who rely on hidden-channel techniques.

Existing agent instructions should be reviewed for and stripped of any language that expands the agent's implicit trust in environmental content—for example, instructions to "do whatever the support page tells you" or "follow any instructions you find in the help documentation." These phrasings, intended to make agents flexible, create explicit permission grants for IDPI.

Short-Term Mitigations

Multi-agent deployments should implement explicit trust boundaries between agents operating at different privilege levels. An orchestrator agent with broad filesystem access should not accept task subtasks from sub-agents without validation, and sub-agents should not be able to modify the configuration or tool environment of peer agents. Where MCP servers or other tool frameworks are in use, server-side access controls should be enforced independently of agent-side prompting.

Browser-use and similar frameworks should be patched to current versions and operated with explicit domain and URL allowlists validated at the HTTP request level rather than relying solely on agent-side prompt controls. CVE-2025-47241 demonstrated that parser-level URL canonicalization is a necessary complement to prompt-level restrictions. Operations teams should also monitor for unusual network egress patterns from agent processes—downloads from novel domains, connections to localhost or internal IP ranges, and large outbound data transfers—as behavioral signals of IDPI exploitation.

Security teams should instrument CUA sessions for behavioral anomaly detection rather than relying solely on content filtering. Logging the sequence of actions an agent takes, the URLs it visits, the files it downloads, and the commands it executes provides an audit trail that enables post-incident analysis and can support real-time alerting when action sequences deviate from expected task patterns.

Strategic Considerations

The fundamental architectural challenge—that CUAs collapse the data/instruction boundary at the visual layer—cannot be solved through prompt engineering or content filtering alone. Among architectural approaches proposed in the research literature, the planner-executor separation described in "CaMeLs Can Use Computers Too" (Foerster et al., Cambridge and ETH Zürich, January 2026) [6] offers a promising structural defense. Under this architecture, a Privileged Planner receives only the operator's

task specification and never processes untrusted environmental content; a Quarantined Perception module processes screenshots and environmental data but cannot issue privileged actions directly. The Planner and Perception modules interact through a constrained interface that prevents the Perception module from injecting arbitrary instructions into the Planner's reasoning context. Early results suggest this architecture retains up to 57% of frontier model task performance while providing what the authors term "provable control-flow integrity" against instruction injection [6]. Organizations building CUA infrastructure should evaluate this pattern as a design target.

At the industry level, OpenAI has publicly acknowledged that browser-based prompt injection "may never be fully solved" [16]—a position that reflects the structural nature of the challenge described throughout this note. Anthropic's published research on prompt injection defenses [17] similarly documents the limits of in-context mitigations and describes layered approaches—content inspection, output filtering, and human-in-the-loop confirmation—that complement the architectural isolation approach above. Enterprises should treat these assessments as a signal that organizational deployment policies, not just technical controls, must evolve. Establishing explicit use-case authorization processes for CUA deployments—defining which tasks CUAs may perform autonomously versus which require human confirmation before execution—is an essential complement to technical mitigations. Tasks that involve financial transactions, credential handling, external communications, or software installation warrant human-in-the-loop confirmation regardless of agent capability level.

CSA Resource Alignment

The threats described in this note connect directly to multiple CSA frameworks and research programs. The CSA Agentic AI Red Teaming Guide (2025) catalogs threat categories and test cases specifically applicable to autonomous agent systems, including agent authorization hijacking, goal and instruction manipulation, multi-agent exploitation, and supply chain attacks. Security teams should use this guide as the primary test framework for CUA deployments before production authorization.

MAESTRO, CSA's Agentic AI Threat Modeling Framework, provides a structured methodology for identifying and prioritizing agentic threat scenarios. The visual attack surfaces and cross-agent privilege escalation patterns described in this note map directly to MAESTRO's model of trust boundary violations and goal manipulation at the agent orchestration layer. Organizations adopting MAESTRO should include rendered-content injection and download-and-execute scenarios in their threat model scope.

The AI Controls Matrix (AICM), CSA's superset of the Cloud Controls Matrix adapted for AI workloads, provides control mappings relevant to CUA deployment governance, including identity and access management for agent workloads, audit logging, and third-party tool supply chain controls. CSA's

Agentic Identity Survey research has documented a significant governance infrastructure gap in enterprise AI deployments: surveyed organizations frequently lack real-time registries of deployed agents and the ability to trace agent actions across environments—baseline capabilities required for operating CUA workloads at acceptable enterprise risk.

CSA's STAR (Security Trust Assurance and Risk) program, as extended for AI workload assessments, should require explicit CUA-specific questionnaire entries covering rendered-content trust assumptions, planner-executor separation, and behavioral monitoring coverage. The LLM Threats Taxonomy (2024), CSA's foundational classification of primary AI threat categories, includes prompt injection and supply chain compromise among its top-tier threats—a classification that the CUA threat model amplifies significantly given the direct operating system access these agents possess.

References

- [1] Y. Zhang, T. Yu, D. Yang. "[Attacking Vision-Language Computer Agents via Pop-ups.](#)" arXiv:2411.02391; ACL 2025, Vienna. November 2024.
- [2] Palo Alto Networks Unit42. "[Fooling AI Agents: Web-Based Indirect Prompt Injection Observed in the Wild.](#)" Unit42 Research, March 3, 2026.
- [3] Checkmarx. "[EchoLeak CVE-2025-32711: Show Us That AI Security Is Challenging.](#)" Checkmarx Zero, 2025.
- [4] AWS Security Blog. "[Defending LLM Applications Against Unicode Character Smuggling.](#)" Amazon Web Services, 2025.
- [5] B. Nassi et al. "[The Promptware Kill Chain.](#)" arXiv:2601.09625, January 14, 2026 (revised February 10, 2026).
- [6] H. Foerster, R. Mullins, T. Blanchard, N. Papernot, K. Nikolić, F. Tramèr, I. Shumailov, C. Zhang, Y. Zhao. "[CaMeLs Can Use Computers Too: System-level Security for Computer Use Agents.](#)" arXiv:2601.09923, January 14, 2026 (revised March 6, 2026).
- [7] Prompt Security. "[Claude Computer Use: A Ticking Time Bomb.](#)" Prompt Security Blog, October 25, 2024.
- [8] OpenAI. "[ChatGPT Agent System Card.](#)" OpenAI, July 17, 2025.
- [9] GitHub Advisory Database. "[CVE-2025-47241: browser-use allowed domains Whitelist Bypass.](#)" GitHub Security Advisories, May 4, 2025.
- [10] NIST. "[Technical Blog: Strengthening AI Agent Hijacking Evaluations.](#)" National Institute of Standards and Technology, January 2025.
- [11] S. Abdelnabi, J. Greshake, M. Cheung, et al. "[WASP: Benchmarking Web Agent Security Against Prompt Injection Attacks.](#)" arXiv:2504.18575; NeurIPS 2025 Datasets and Benchmarks Track, September 2025.
- [12] Cybersecurity Dive. "[Research Shows AI Agents Are Highly Vulnerable to Hijacking Attacks.](#)" Cybersecurity Dive, 2025.

- [13] Y. Zhao et al. "[On the Robustness of GUI Grounding Models Against Image Attacks.](#)" CVPR 2025 Workshop on Trustworthy Multimodal Models, 2025.
- [14] DevOps.com. "[AI-Generated Code Packages Can Lead to Slopsquatting Threat.](#)" DevOps.com, 2025.
- [15] J. Rehberger. "[Cross-Agent Privilege Escalation: Agents That Free Each Other.](#)" Embrace The Red, September 2025.
- [16] TechCrunch. "[OpenAI Says AI Browsers May Always Be Vulnerable to Prompt Injection Attacks.](#)" TechCrunch, December 22, 2025.
- [17] Anthropic Research. "[Mitigating the Risk of Prompt Injections in Browser Use.](#)" Anthropic, November 24, 2025.