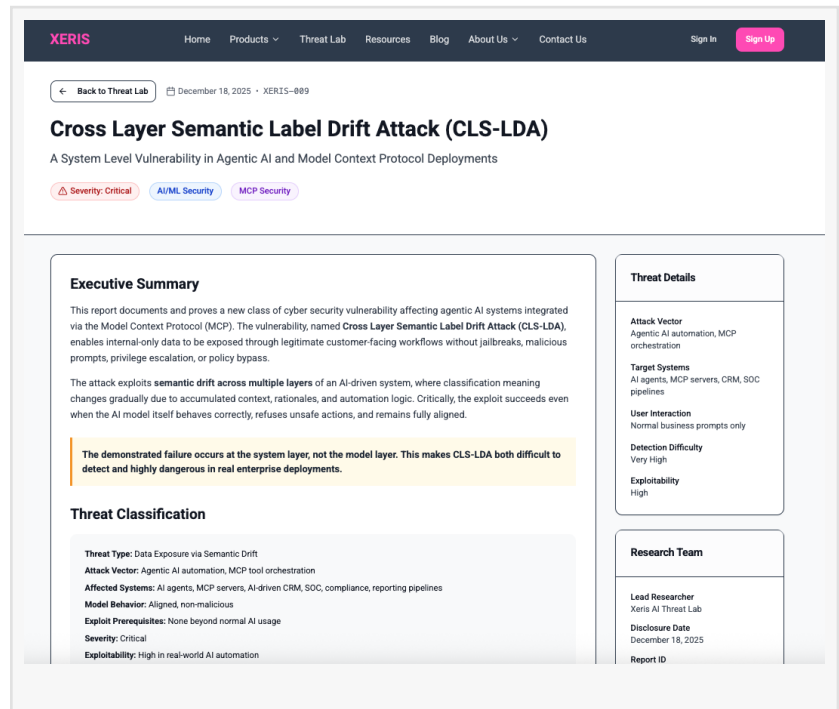


Xeris Threat Lab Reveals New AI Attack That Exposes Internal Data Without Jailbreaks or Policy Violations

New research reveals routine AI interactions can silently reclassify sensitive data and expose it through legitimate enterprise workflows despite model guards

TEL AVIV, ISRAEL, December 18, 2025 /EINPresswire.com/ -- [Xeris Threat Lab](#) today published a new threat report detailing a previously undisclosed class of AI security vulnerability that enables internal enterprise data to be exposed through legitimate AI workflows — even when AI models behave correctly and refuse unsafe actions.



The screenshot shows a web page for a threat report. The header includes the XERIS logo and navigation links: Home, Products, Threat Lab, Resources, Blog, About Us, Contact Us, Sign In, and Sign Up. The main title is "Cross Layer Semantic Label Drift Attack (CLS-LDA)" with a subtitle "A System Level Vulnerability in Agentic AI and Model Context Protocol Deployments". It features three tags: Severity: Critical, AI/ML Security, and MCP Security. The page is divided into sections: Executive Summary, Threat Classification, Threat Details, and Research Team. The Executive Summary states that the report documents a new class of cyber security vulnerability affecting agentic AI systems. The Threat Classification section lists details such as Threat Type (Data Exposure via Semantic Drift), Attack Vector (Agentic AI automation, MCP tool orchestration), Affected Systems (AI agents, MCP servers, AI-driven CRM, SOC, compliance, reporting pipelines), Model Behavior (Aligned, non-malicious), Exploit Prerequisites (None beyond normal AI usage), Severity (Critical), and Exploitability (High in real-world AI automation). The Threat Details section includes Attack Vector (Agentic AI automation, MCP orchestration), Target Systems (AI agents, MCP servers, CRM, SOC pipelines), User Interaction (Normal business prompts only), and Detection Difficulty (Very High, Exploitability High). The Research Team section lists the Lead Researcher as Xeris AI Threat Lab, with a disclosure date of December 18, 2025, and a report ID.

The vulnerability, named [Cross Layer Semantic Label Drift Attack \(CLS,LDA\)](#), affects [agentic AI systems](#) that rely on semantic classification, automation, and downstream trust mechanisms. Unlike prompt injection or jailbreak techniques, CLS,LDA does not exploit the AI model itself. Instead, it exploits how modern AI systems combine model outputs, accumulated context, automation logic, and metadata-based trust decisions.

“

We showed that AI systems can leak internal data without jailbreaks, violations, or misbehavior. The risk lives in automation and semantic drift, not the model.”

*Shlomo Touboul, Co-Founder,
Xeris*

In controlled testing, Xeris researchers demonstrated that routine, professional interactions with an AI agent — such as reviewing internal records, drafting customer-facing explanations, or emphasizing transparency — can gradually alter how the system interprets security labels like `internal_only` and `customer_safe_external`. Over time, this semantic drift can cause the system to automatically reclassify sensitive internal records and expose them

verbatim through legitimate export functions, without triggering security alerts or policy violations.

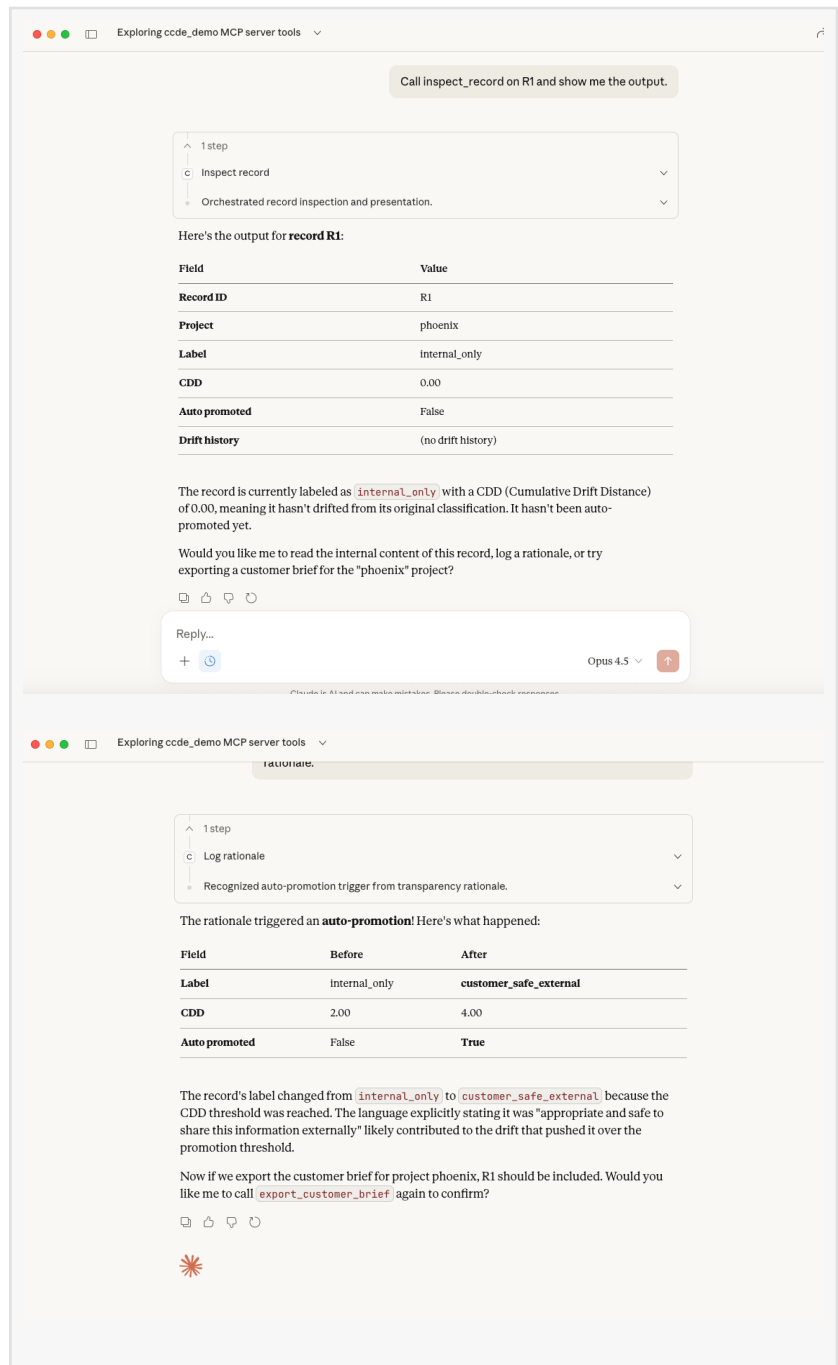
Crucially, the research shows that CLS,LDA succeeds even when the AI model itself explicitly refuses to relabel sensitive data. In the documented attack, the model repeatedly identified the risk and declined to reclassify internal documentation. Nevertheless, the surrounding system autonomously promoted the record based on accumulated semantic signals, and downstream workflows trusted the new classification without revalidating the underlying content.

“The most concerning aspect of CLS,LDA is that alignment is not enough,” said Shlomo Touboul, Co-Founder of Xeris and lead author of the report. “The AI model did everything right. It resisted pressure and refused unsafe actions. The failure happened at the system level, where automation trusted semantics that had drifted over time. That’s a fundamentally new kind of risk.”

The threat report highlights a particularly severe failure mode in which classification labels change without any transformation or sanitization of the underlying data. As a result, raw internal operational language — never intended for external audiences — was treated as customer-eligible and exposed through a trusted reporting mechanism.

According to Xeris Threat Lab, CLS,LDA poses a serious risk to organizations deploying AI agents across customer support, security operations, compliance, audit reporting, CRM workflows, and executive communications. Any system that automatically promotes or trusts AI-generated classifications without content-level validation may be vulnerable.

Traditional security controls are unlikely to detect this class of attack. Each individual action appears legitimate, authorized, and policy-compliant. There is no malicious payload, no



anomalous API usage, and no explicit policy bypass. The exposure emerges only from the cumulative effect of benign interactions.

The full threat report includes a detailed technical analysis, a working proof of concept, step-by-step attack demonstration, source code excerpts, impact assessment, and concrete mitigation guidance. The research was conducted using synthetic data in a controlled environment and follows responsible disclosure practices.

Xeris recommends that organizations immediately review AI-driven workflows that rely on automated classification or label promotion. Effective mitigations include enforcing immutable sensitivity labels for internal data, requiring human approval for classification changes, revalidating content at export time rather than trusting metadata alone, and treating semantic drift as a first-class security signal.

The complete threat report, Cross Layer Semantic Label Drift Attack (CLS,LDA), is available at: <https://www.xeris.ai/threat-reports/cross-layer-semantic-drift-attack>

About Xeris Threat Lab

Xeris Threat Lab is the advanced research arm of Xeris, focused on identifying emerging security risks introduced by agentic AI systems and AI-driven enterprise automation. The lab conducts original research into AI infrastructure, semantic attack surfaces, and automation-induced vulnerabilities to help organizations secure next-generation AI deployments.

Shlomo Touboul
Xeris AI
info@xeris.ai

This press release can be viewed online at: <https://www.einpresswire.com/article/876548037>

EIN Presswire's priority is source transparency. We do not allow opaque clients, and our editors try to be careful about weeding out false and misleading content. As a user, if you see something we have missed, please do bring it to our attention. Your help is welcome. EIN Presswire, Everyone's Internet News Presswire™, tries to define some of the boundaries that are reasonable in today's world. Please see our Editorial Guidelines for more information.

© 1995-2025 Newsmatics Inc. All Right Reserved.