

Zero Latency Threat Detection

Gurucharan Raghunathan^{1*}, Michael Wilson Rebello¹

¹Trench Security, Bangalore, India

¹Email: guru@trenchsecurity.ai; michael@trenchsecurity.ai

*Corresponding Author: guru@trenchsecurity.ai

ARTICLE INFO

Received: 15 Jan 2026

Revised: 26 Feb 2026

Accepted: 10 March 2026

ABSTRACT

In the age of AI-powered threats, the physics of cyber warfare has changed. While defensive strategies have spent the last decade perfecting data aggregation, adversaries have perfected speed. Today, automated attacks execute in milliseconds, yet our industry's standard detection processes are measured in minutes, hours, or weeks. This temporal disconnect, the "Latency Gap" is where the modern breach lives. For too long, security analytics has been perceived as a Data Problem and we were convinced to ingest everything and ended up with a data overload bottleneck. We hoard petabytes of logs to find the "needle in the haystack." This model is obsolete. AI has handed adversaries a new, defining moat: Velocity. AI-native threats do not pause for human cognition. Relying on legacy, data-centric detection models to fight these threats is akin to fighting drone warfare with ground patrol. To combat this, we must shift our mindset from Data to Time. You can have the most accurate threat detection rule in existence, but if it triggers 48 hours or even 48 seconds after the event, it is not a defense; it is forensics. The collateral damage has already occurred; the attacker is far ahead in the kill chain. Just as the Zero Trust framework revolutionized security by shifting our focus from the "Perimeter" to "Identity" (Data), we must now evolve further. We introduce the Zero Latency Threat Detection (ZLTD) framework. This architecture accepts that in a world of instantaneous execution, the only effective defense is instantaneous detection. We must stop hunting for the needle and start catching it before it lands. Cybersecurity Mesh Architecture (CSMA) is the right approach to apply ZLTD that moves detection to the edge, reducing SIEM data gravity costs while delivering real-time, high-fidelity threat detection.

Keywords: Artificial Intelligence (AI), Cybersecurity Mesh Architecture (CSMA), Latency Gap, Real-Time Threat Detection, SIEM Optimization, Zero Latency Threat Detection (ZLTD)

1. Introduction

The acceleration of AI-driven attack automation has fundamentally altered the temporal dynamics of cybersecurity, compressing exploitation cycles from hours to milliseconds. Traditional detection architectures, designed around centralized log aggregation and retrospective analysis, are increasingly misaligned with this new threat velocity. This widening latency gap enables adversaries to advance through the kill chain before defensive controls can respond. To address this structural imbalance, this paper introduces the Zero Latency Threat Detection (ZLTD) framework, a purpose-driven architecture that shifts security from data-centric processing to time-centric execution through decentralized, edge-based intelligence.

A. The Velocity Gap: Attack Speed vs. Detection Capability

The traditional barrier to entry for cyber adversaries has collapsed. AI has democratized the attacker by augmenting both traditional APTs and lower-skill adversaries, shifting the landscape from "Advanced Persistent Threats" (APTs) to "Automated Persistent Threats."

B. The Evidence: An Exponential Zero Day Explosion

We are losing the war because the opportunity for attacks (Vulnerabilities) and the volume of threat intelligence (Artifacts) are diverging faster than human teams can handle.

1) *The Vulnerability Flood (~140/Day)*: In 2018, teams triaged ~45 new Zero day CVEs/day. By Nov 2025, that number has tripled to ~140/day. Every 10 minutes, a new door opens on the global attack surface.

2) *The Intelligence Avalanche (~42.5 Million/Day)*: According to the *Microsoft Digital Defense Report 2025*, the industry now faces 4.5 million new malware files and 38 million identity risk detections daily.

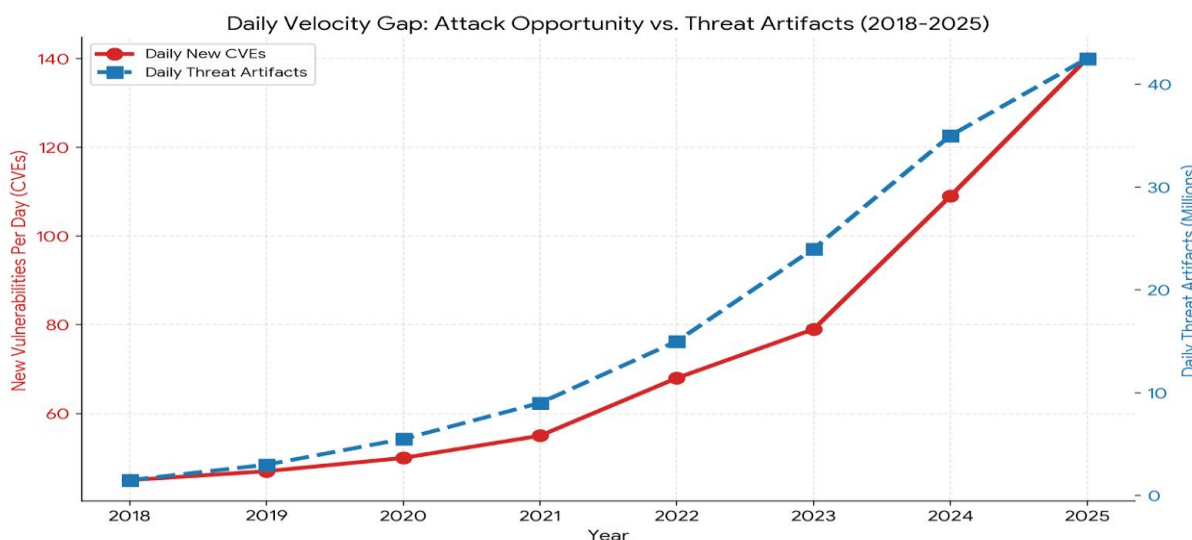


Fig. 1 Exponential growth in daily CVEs and threat artifacts (2018–2025) highlights the widening velocity gap between attack opportunity and detection capacity.

C. The Root Cause: Why Velocity is Exploding

We are not just facing "more" attacks; we are facing a fundamental capability upgrade. The "Latency Gap" exists because the barrier to entry has collapsed due to three specific market forces:

- 1) *Attack Commoditization (RaaS)*: Cybercrime has adopted the SaaS model. Ransomware-as-a-Service decouples "Malware Authors" from "Operators," allowing low-skill affiliates to rent enterprise-grade tools without needing technical expertise.
- 2) *The "Skill Floor" Collapse*: Research from the UC Berkeley Risk and Security Lab indicates that Generative AI has lowered the entry barrier, allowing attackers to automate complex tasks, like coding variants and phishing, turning attacks that cost days of labor into pennies of compute.
- 3) *Infrastructure as a Utility*: As detailed by OPSWAT, high-end infrastructure is no longer the domain of Nation States. Attackers now utilize AI agents to automate the entire kill chain, executing reconnaissance and payload generation at machine speed that outpaces human response teams.

D. Case Study: How AI can automate E2E attack from a CVE

The AI-Automated SSRF Attack The real danger is not just the number of vulnerabilities, but the speed of exploitation. Consider a standard Server-Side Request Forgery (SSRF) vulnerability:

- 1) *Legacy Attack*: A human actor identifies the SSRF, manually probes internal ports, and spends hours crafting a payload to bypass filters.
- 2) *AI-Driven Attack*: An AI agent ingests the new CVE, automates the end-to-end attack path instantly. It identifies the SSRF vulnerability, auto-scans for internal targets, chains it with a secondary privilege escalation exploit, and crafts a bespoke payload to exfiltrate cloud credentials—all in milliseconds.

This automation turns what used to be a "Low Severity" misconfiguration into a "Critical" breach faster than any human analyst can triage the ticket.

E. The Speed of Loss: Mean Time to Exfil (MTTE)

The "Velocity Gap" is best illustrated by the collapse of "Dwell Time." According to Palo Alto Networks (2025), the window for defense has effectively closed:

- 1) *100x Acceleration*: Ransomware attacks that took 9 days in 2021 now complete in as little as 25 minutes.
- 2) *The AI Multiplier*: This speed is driven by AI-powered "Multi-Extortion" campaigns. 82.6% of phishing emails now utilize AI to automate the deployment of quadruple extortion tactics (Encryption, Theft, DoS, and Harassment).
- 3) *The Reality*: If the Mean Time to Exfil (MTTE) is 25 minutes, a detection process measured in hours is not security, it is merely observation of a breach in progress.

F. The Three Velocities of Failure

- 1) *Asset Velocity (Visibility Collapse)*: Assets created by cloud, SaaS, and Shadow AI appear, process data, and disappear in minutes, faster than traditional discovery and scan cycles can track.
- 2) *TTP Velocity (AI-Driven Mutation)*: Adversary techniques are no longer static; AI enables real-time mutation to evade fixed, rule-based detection.
- 3) *Intelligence Velocity (Correlation Overload)*: Threat intelligence arrives at machine scale, making real-time correlation across assets, behaviors, and signals infeasible for human-led or query-based systems.

The current model of Static Detection Rules fails because it attempts to manually correlate these exploding variables. We cannot fight the velocity of AI with the latency of a query language.

2. The Solution: Introducing 'High Recall' Detection

To achieve *Zero Latency Threat Detection (ZLTD)*, we must fundamentally alter the engineering mindset. The legacy model is constrained by human bandwidth; the ZLTD model releases those constraints using the power of *AI Agents*.

G. The Paradigm Shift: From "Precision-First" to "Recall-First"

The fundamental bottleneck in traditional "Detection-as-Code" is the Human Detection Mindset. Because human analysts cannot process 10,000 alerts a day, detection engineers are forced to tune rules for High Precision (Low False Positives). With AI Agents, we can flip the model. We shift the focus to High Recall (Catch everything).

- 1) *Legacy Model (Precision)*: Human expertise plays a critical role in creating and fine tuning of rules for precision in accordance with the ever evolving threat landscape. Result: We miss detecting new age

attacks. Because human analysts cannot process 10,000 alerts/day, engineers tune rules to be highly specific. Result: We miss slight variations of attacks (Low Recall) despite Defense in Depth architecture. 2) *Zero Latency Model (High Recall)*: AI Agents have infinite bandwidth. We can create and tune rules for High Recall (Catch everything). The Agent acts as the first tier of analysis, instantly "Red Teaming" every signal to filter noise, leaving humans with only confirmed threats.

The Operational Shift: Overcoming the Precision Paradox

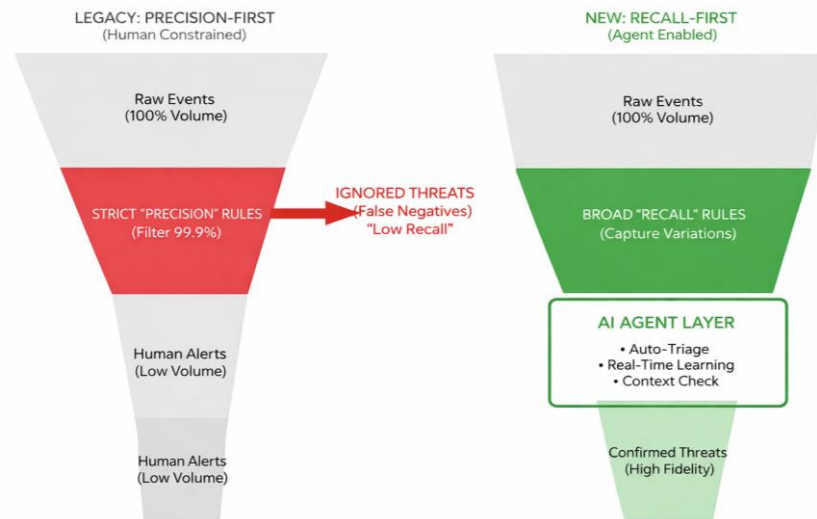


Fig.2 Transition from precision-first, human-constrained detection to recall-first, agent-enabled architecture resolves the precision paradox by combining broad signal capture with AI-driven validation for high-fidelity threat confirmation.

The AI Agent acts as the first tier of analysis, capable of investigating thousands of "potential" signals instantly to filter out the noise, leaving the human with only high-fidelity and confirmed threats. This model enables the security teams to achieve dynamic, accurate and real-time high fidelity detection rules to protect their critical assets. Ideal model should let human experts focus on *Precision* and AI systems to meet *Recall* at scale.

H. The Proof: Kill Chain Scenario (Phishing)

Kill Chain Timeline Comparison: Legacy vs. High-Recall AI Agent

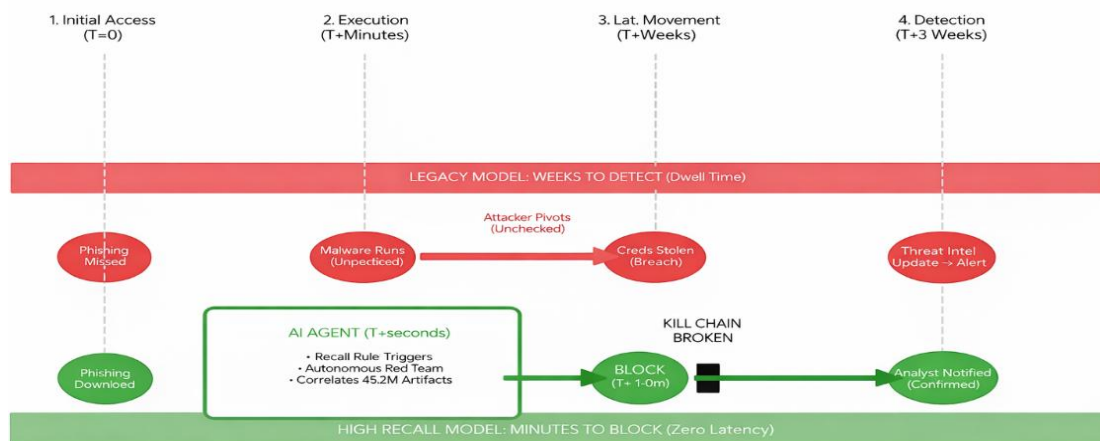


Fig.3 Scenario: A user downloads a malicious .zip file from a phishing email.

- 1) *Legacy Path (Red)*: Strict rules miss the new file hash. The attack proceeds for *weeks* until a retrospective Threat Intel update triggers an alert. Credentials are already stolen. The current manual process gap in Defense in Depth also delays the existing model to capture the lateral movement signals.
- 2) *High-Recall Path (Green)*: The ZLTD framework flags the unusual download in *minutes*. The AI Agent autonomously sandboxes the file, correlates it with global artifacts, and blocks the attack in <20 minutes, breaking the kill chain before lateral movement occurs.

3. Strategic Architecture: The AI-Powered CSMA

To implement *Zero Latency Threat Detection (ZLTD)*, we operationalize *Gartner’s Cybersecurity Mesh Architecture (CSMA)*. With rapid cloud adoption in the past decades, the industry’s answer to complexity has been *Centralization*. We convinced ourselves that if we just funneled every log, signal, and event into a massive SIEM or Data Lake, we could find the threat. This approach created a massive "Data Gravity" problem. We are spending more budget on *moving and storing* data than on *securing* it.

Additionally, establishing data pipelines from multiple applications impacts the end user application performance and adds overhead on network bandwidth. With the ZLTD framework, instead of building hard-coded pipes between tools, we use *AI Agents* as the dynamic "connective tissue" that enables decentralized enforcement.

I. The "Silo" Advantage

Humans struggle with silos; AI thrives in them.

- 1) *Legacy Workflow (Centralized SIEM)*: We spend budget shipping logs to a central lake to analyze them. This creates a massive "Ingestion Tax" and latency lag. Additionally, when opting for data lake environments, there is a significant risk of Data Security and Privacy since the identity controls at the source are often disrupted at the destination.
- 2) *Zero Latency Workflow (Decentralized Mesh)*: Data stays put. The AI Agent connects to the Endpoint (Silo A) and Cloud (Silo B), correlating signals in real-time without moving the data.

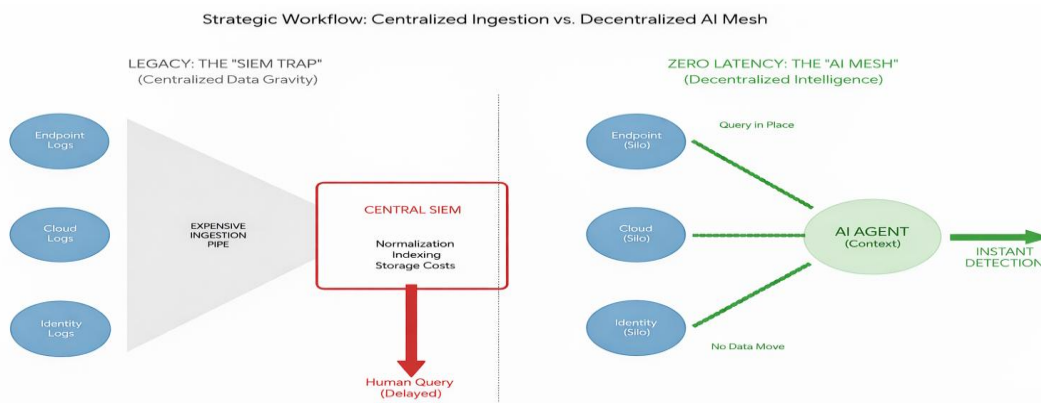


Fig.4 Comparison of centralized SIEM ingestion versus decentralized AI mesh architecture illustrates how edge-based intelligence eliminates data gravity and enables instant, in-place threat detection.

J. Disrupting the SIEM Layer

This framework redefines the SIEM:

- 1) *From*: The "Brain" of Detection (High Cost/High Latency).

2) *To:* The "System of Record" for Compliance (Low Cost/High Storage).

By moving detection to the Mesh, organizations unlock speed while drastically reducing data centralization costs. Organizations can stop paying premium "Hot Storage" rates for detection data. They can keep their detection logic at the edge (where the AI is) and only send *confirmed alerts* and *compliance logs* to the SIEM. This yields a massive cost reduction while simultaneously increasing detection speed. Fundamentally, this framework drives us to reimagine what SIEM can do in an AI era.

K. *Technological Enabler: The Shift to Real-Time Data Primitives*

To match the velocity of AI-driven threats, we must abandon legacy, general-purpose Data Warehouses in favor of specialized *Real-Time OLAP* architectures. By adopting *Open Table Formats* and high-performance columnar indexing, we decouple compute from storage, eliminating the "ingestion tax" and proprietary lock-in of traditional SIEMs. This shift transforms detection from a batch process measured in minutes to a continuous query measured in milliseconds, enabling instant search across petabytes of telemetry at a fraction of the cost.

4. The Business Impact: Returning to the 'WHY'

Applying Simon Sinek’s Golden Circle (Start with Why) framework, we can visualize the strategic pivot required. The legacy model is broken because it starts with the wrong question for the information security leaders. Zero Latency Threat Detection corrects the course by starting with the fundamental purpose of security which is detecting threats instantly and protecting the business.

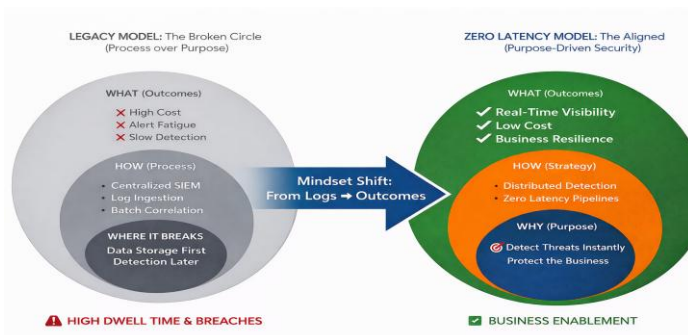


Fig.5 Evolution from a process-driven, centralized security model to a purpose-driven Zero Latency architecture aligns strategy, detection, and business outcomes to reduce dwell time and enable operational resilience.

L. *The Strategic Pivot*

1) *Legacy (Left):* We focused on WHAT (Logs/Compliance) and built a HOW (Centralized SIEM) that became a bottleneck. The WHY (Protection) was lost.

2) *Zero Latency (Right):* We start with WHY (Detect Threats & Protect Revenue). This dictates the HOW (Distributed ZLTD & CSMA). The result is WHAT businesses actually need: Visibility, Velocity, and Cost Efficiency.

M. *CISO Success Metrics: The Business Impact of Zero Latency*

1) *Cost Efficiency:* Delivers 40-60% reduction in SIEM costs by eliminating the "Ingestion Tax", paying only for high-fidelity signals, not noise.

2) *Risk Exposure:* Collapses *Dwell Time* from weeks to minutes, mathematically eliminating the window for data exfiltration.

3) *Productivity*: Acts as a *10x Force Multiplier* by automating noise filtering and tuning, allowing existing teams to handle AI-scale threat volume.

N. The Regulatory Stopwatch: Speed is Now Law

Global regulators have recognized the "Velocity Gap" and responded with strict reporting mandates. The era of "weeks to investigate" is over; organizations are now legally required to detect, scope, and report material incidents within hours.

Legacy batch-detection architectures (with 24+ hour latency) effectively guarantee non-compliance with these new standards:

SI.NO	Regulation	Jurisdiction	The "Speed Limit" (Mandate)	The Risk
I	SEC Rules	USA (Public Corp)	4 Business Days to determine materiality and disclose.	Detection delays consume the investigation window, forcing premature (or inaccurate) 8-K filings.
II	GDPR	EU / Global	72 Hours to notify supervisory authorities of a breach.	Fines up to 4% of global revenue for failure to report "without undue delay."
III	CIRCA	USA (Critical Infra)	72 Hours for covered incidents; 24 Hours for ransomware payments.	CISA mandates rapid reporting to prevent systemic contagion.
IV	NIS2 / DORA	EU (Critical/Finance)	24 Hours ("Early Warning") for significant incidents.	Financial entities face penalties for failing to provide "near real-time" situational awareness.
V	DPDP Act 2023	India	Without Undue Delay (< 72 Hours implied) to notify Board & Users.	Penalties up to ₹250 Crore for failure to take reasonable safeguards or notify breaches.

5. Conclusion

For fifteen years, we attempted to secure the enterprise by building bigger data lakes, only to drown in costs while adversaries moved faster. *Zero Latency Threat Detection* corrects this fundamental architectural error. By shifting from a centralised storage-centric legacy model to an intelligence-centric decentralized Mesh, organizations achieve the "Holy Grail" of cyber defense: *Lower Costs, Higher Velocity, and Total Visibility*.

We are no longer just buying tools; we are buying *time*. In a landscape where AI attacks execute in milliseconds, the ability to detect and respond in the same clock cycle is not a luxury, it is the baseline for business resilience. The technology is proven. The ROI is validated. The only decision left is to stop funding the latency gap and start closing it.

6. Acknowledgement

we sincerely thank the following subject matter experts for their technical validation, domain insights, and strategic guidance in shaping the Zero Latency Threat Detection (ZLTD) framework:

- 1) *Subhro Banerjee*, Deputy CISO at a global pharmaceutical organization;
- 2) *Sammit Potdar*, former Global CISO across pharmaceutical, oil and gas sectors;
- 3) *Senthil Kumar Iyyappan*, CISO of a U.S.-based fintech firm;
- 4) *Arindam Ghose*, Global CISO Advisor in digital operations;
- 5) *Dr. Ganesh Nagaraj*, Global InfoSec Advisor and SOC researcher;
- 6) *Lokesh Kannaiyan*, Cybersecurity Product Management Specialist;
- 7) *Aravindh Seenevasan*, Incident Response and SOC Expert.

7. References

- [1] NIST National Vulnerability Database (NVD) Dashboard: <https://nvd.nist.gov/general/visualizations/vulnerability-visualizations>
- [2] Gartner Cybersecurity Mesh Architecture (CSMA): <https://www.gartner.com/en/information-technology/glossary/cybersecurity-mesh>
- [3] Microsoft Digital Defense Report 2025: <https://www.microsoft.com/en-us/security/security-insider/microsoft-digital-defense-report-2025>
- [4] Generative AI & LLMs are easily available from Berkeley- <https://brsl.berkeley.edu/cybercrime-gets-an-upgrade/> because
- [5] Sophisticated LLMs & infra readily available - <https://www.opswat.com/blog/ai-hacking-how-hackers-use-artificial-intelligence-in-cyberattacks>
- [6] IBM X-Force Threat Intelligence Index 2025: <https://www.ibm.com/reports/threat-intelligence>
- [7] AV-TEST Institute Malware Statistics: <https://www.av-test.org/en/statistics/malware/>
- [8] Palo Alto Networks: The Ransomware Speed Crisis: Why AI Changed the Clock <https://www.paloaltonetworks.com/blog/2025/09/ransomware-speed-crisis/>
- [9] Gartner Continuous Threat Exposure Management (CTEM): <https://www.gartner.com/en/articles/what-is-ctem-and-why-is-it-important>
- [10] NIST Special Publication 800-207 (Zero Trust Architecture): <https://csrc.nist.gov/pubs/sp/800/207/final> Ponemon Institute: The Economics of Security Operations
- [11] (<https://www.scribd.com/document/486285050/Ponemon-Report-The-Economics-of-Security-Operations-Centers>)
- [12] IDC Global Datasphere Forecast (Data Volume Projection): <https://my.idc.com/getdoc.jsp?containerId=US53383425> MITRE ATT&CK Framework: <https://attack.mitre.org/>
- [13] MITRE ATLAS (Adversarial Threat Landscape for AI-Systems): <https://atlas.mitre.org/>

- [14] Enterprise SIEMs are failing to detect 79% of all MITRE ATT&CK techniques: <https://cardinalops.com/white-papers/2025-state-of-siem-report-download/>.
- [15] Cost of a Data Breach Report 2025: <https://www.ibm.com/reports/data-breach>
- [16] M-Trends 2025: <https://cloud.google.com/security/resources/m-trends>
- [17] <https://market.us/report/cybersecurity-analytics-market/>
- [18] Simon Sinek's Golden Circle: <https://simonsinek.com/golden-circle/>
- [19] CISA Secure by Design Guidelines: <https://www.cisa.gov/securebydesign>
- [20] A comparative study of Delta Parquet, Iceberg, and Hudi for automotive data engineering use cases. International Journal of Computer
- [21] Science and Engineering, 12(17), 104. <https://doi.org/10.14445/23488387/IJCSE-V12I17P104>
- [22] 2016 IEEE International Conference on Cluster Computing (CLUSTER), 354–363. <https://doi.org/10.1109/cluster.2016.29>