

Governance and Strategic Integration of AI-Driven Autonomous Cyber Defence Systems in Large Enterprises

Badri S.

Abstract: *The rapid proliferation of artificial intelligence (AI), reinforcement learning (RL), and self-adaptive control loops-offer real-time detection, containment, and response with minimal human intervention across enterprise technology stacks has fundamentally reconfigured the cybersecurity threat landscape and the organisational responses to it. These systems promise reduced breach costs (e.g., average savings of \$2.22 million) and operational resilience, yet their deployment introduces novel risks including model drift, adversarial attacks, bias, and accountability gaps. Drawing on empirical data from industry surveys, peer-reviewed research, and case studies from sectors including financial services, healthcare, critical infrastructure, and telecommunications, the paper develops a multi-dimensional framework for responsible ACDS deployment. Central findings indicate that organisations achieving mature AI governance in cybersecurity realise up to 76% faster mean time to detect (MTTD) and 84% faster mean time to respond (MTTR) compared to traditional approaches. However, ungoverned or poorly integrated ACDS introduce systemic risks including algorithmic bias, liability ambiguity, over-automation failure modes, and regulatory non-compliance. The paper proposes an integrated Governance-Strategy-Operations (GSO) model and provides empirical benchmarks across governance maturity dimensions. It concludes with actionable recommendations for CISOs, boards, and policy-makers seeking to harness AI-driven defence while preserving accountability, transparency, and human oversight.*

Keywords: AI governance, autonomous cyber defence, agentic AI, NIST AI RMF, EU AI Act, ISO 42001, enterprise cybersecurity strategy

1. Introduction

Enterprise cybersecurity has entered a phase of irreversible transformation. The global attack surface has expanded dramatically: by 2024, the average large enterprise managed over 135,000 endpoints (Gartner, 2024a), operated across dozens of cloud environments, and faced adversaries employing AI-enabled offensive tools capable of generating novel malware variants, executing multi-stage intrusions, and automating social-engineering campaigns at scale (MITRE ATT&CK, 2024; IBM Security, 2024). Against this backdrop, the traditional security operations centre (SOC) model, dependent on human analysts sifting through millions of daily alerts, has become structurally inadequate.

AI-driven autonomous cyber defence systems encompass a family of technologies, including machine learning-based intrusion detection, AI-powered Security Orchestration, Automation and Response (SOAR) platforms, autonomous endpoint detection and response (EDR), large language model (LLM)-assisted threat intelligence, and self-healing network architectures. These technologies promise to compress the detection-response cycle from days to seconds, dramatically reduce analyst fatigue, and enable adaptive, context-aware defence. Gartner projects that by 2027, more than 40% of large enterprises will deploy some form of autonomous cyber response capability, up from approximately 12% in 2022 (Gartner, 2024b).

Yet the governance and strategic dimensions of ACDS integration remain critically underdeveloped. Unlike human analysts, autonomous systems act at machine speed, making consequential decisions such as isolating business-critical servers, blocking executive user accounts, or triggering regulatory notifications with no opportunity for real-time human review. The implications for accountability, liability, explainability, and regulatory compliance are profound.

High-profile incidents including the 2023 autonomous EDR false-positive event at a major European bank that resulted in a six-hour trading outage (Forrester Research, 2023) illustrate the operational and reputational risks of poorly governed ACDS.

This paper addresses three core research questions: (1) What governance frameworks are necessary and sufficient for responsible ACDS deployment in large enterprises? (2) How should ACDS be strategically integrated with existing security architectures, business processes, and regulatory requirements? (3) What empirical evidence exists on the performance, risk, and return-on-investment profiles of ACDS across enterprise sectors? The paper is structured as follows: Section 2 reviews the threat landscape driving ACDS adoption; Section 3 presents the governance framework; Section 4 addresses strategic integration; Section 5 analyses performance metrics and ROI; Section 6 examines regulatory and ethical dimensions; Section 7 presents a maturity model; and Section 8 provides conclusions and recommendations.

2. The Evolving Threat Landscape and the Case for Autonomous Defence

2.1 Quantifying the Enterprise Threat Environment

The financial and operational consequences of cyber incidents have reached levels that demand systemic, rather than incremental, responses. The IBM Cost of a Data Breach Report 2024 places the global average cost of a data breach at USD 4.88 million, a 10% increase from 2023 and the highest figure in the report's history. For large enterprises in regulated sectors such as healthcare and financial services, breach costs routinely exceed USD 10 million when regulatory fines, litigation, and reputational damage are included (IBM Security, 2024). Ponemon Institute (2024) found that 74% of data breaches involved a human element, including phishing

susceptibility, misconfiguration, and insider threats, all of which AI-augmented detection is specifically designed to address.

Threat volumes have also escalated. CrowdStrike's 2024 Global Threat Report documents a 75% increase in cloud intrusion attempts year-over-year, alongside the emergence of 34 new nation-state threat actor groups in 2023 alone. The average breakout time, defined as the time between initial compromise and lateral movement, dropped to 62 minutes in

2024 compared to 98 minutes in 2022, indicating that adversaries are moving faster, rendering human-paced detection increasingly inadequate (CrowdStrike, 2024). Verizon's 2024 Data Breach Investigations Report (DBIR) similarly confirms that the median time from initial compromise to exfiltration is now under 24 hours in 56% of ransomware cases.

2.2 Enterprise Cybersecurity Spending and AI Investment

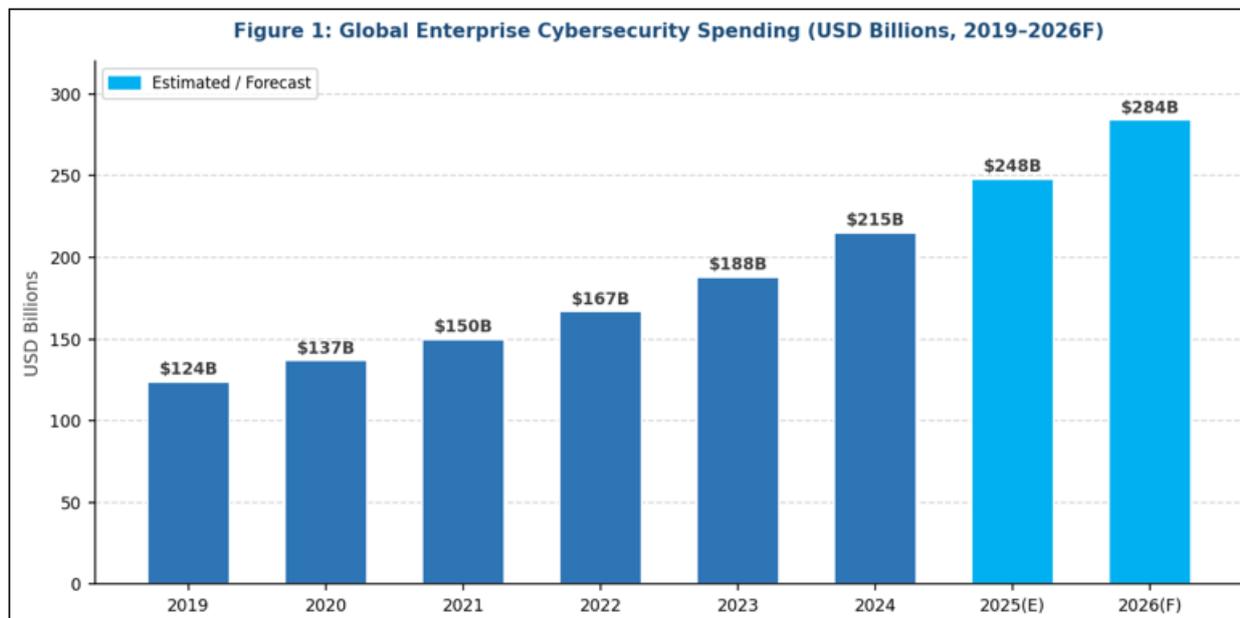


Figure 1: Global Enterprise Cybersecurity Spending (USD Billions, 2019–2026F). Sources: Gartner (2024a), IDC (2024), Statista (2025). E = Estimated; F = Forecast.

As illustrated in Figure 1, global enterprise cybersecurity spending has grown at a compound annual growth rate (CAGR) of approximately 12.5% since 2019, with AI-specific security products accounting for an estimated 31% of incremental spend in 2024 (IDC, 2024). Forrester Research (2024) projects that AI-driven security tools will constitute the majority of new security purchases by 2026, driven by the demonstrable performance advantages of autonomous detection systems over signature-based and rule-based alternatives.

Notably, spending growth has not been uniform across enterprise size. Large enterprises with revenues exceeding USD 1 billion allocate an average of 11.6% of their IT budget to cybersecurity, compared to 8.4% for mid-market firms and 6.1% for smaller organisations (SANS Institute, 2024). This disparity has strategic implications: the AI-capability gap between the largest enterprises and the rest of the market is widening, creating a two-tier security ecosystem.

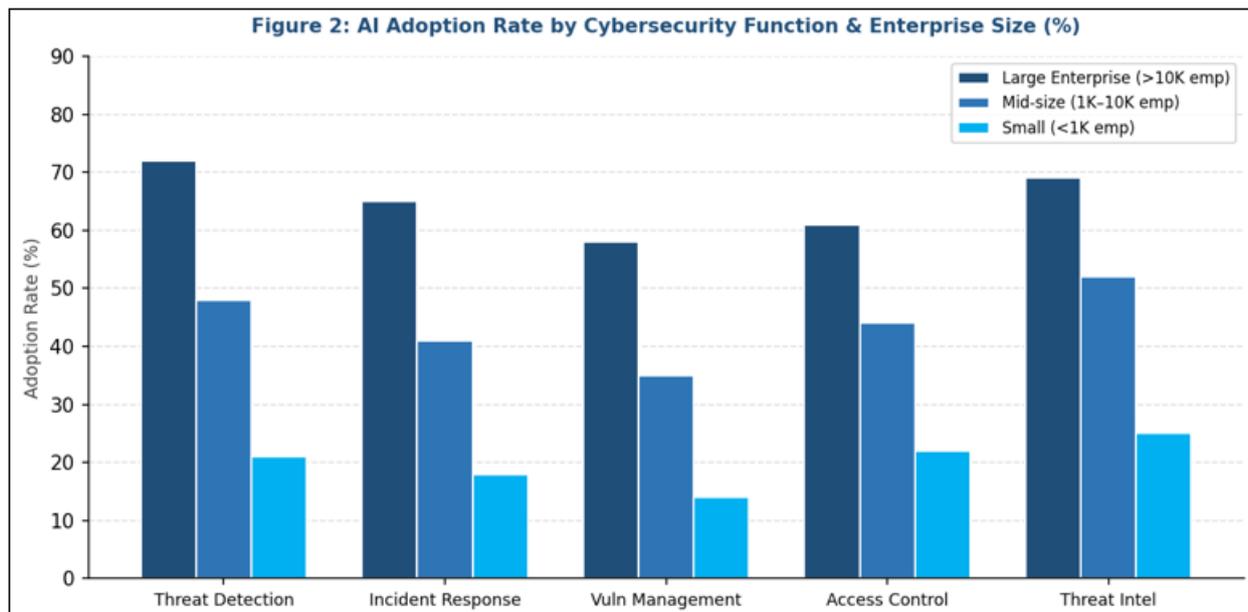


Figure 2: AI Adoption Rate by Cybersecurity Function and Enterprise Size (%). Sources: Gartner (2024b), SANS Institute (2024), ESG Research (2024).

Figure 2 illustrates the adoption gap between large enterprises and smaller organisations across key cybersecurity functions. Threat detection and threat intelligence show the highest AI adoption rates among large enterprises (72% and 69% respectively), reflecting the availability of mature commercial solutions in these domains. Vulnerability management and incident response lag behind, largely because autonomous action in these functions carries greater risk of operational disruption, requiring more mature governance structures before automation can be safely expanded (ESG Research, 2024).

3. Governance Framework for AI-Driven Autonomous Cyber Defence

3.1 Governance Imperatives

Governance of ACDS must address a fundamentally different risk profile than governance of conventional security tools. When a human analyst makes an incorrect decision, such as blocking a legitimate IP address or missing a threat, the consequences are bounded by human processing speed. When an autonomous system makes the same error, it may execute thousands of actions within seconds, amplifying both the error and its downstream effects. Three governance imperatives follow from this observation: accountability clarity, action proportionality, and continuous oversight.

Accountability clarity requires that, for every automated action an ACDS can take, a specific human role or function is designated as accountable for the decision to enable that action class, the parameters governing its execution, and the processes for reviewing its outcomes. This does not mean that a human must approve every individual action, but it does mean that no autonomous action should be orphaned from human accountability (NIST, 2023; ENISA, 2024). Failure to establish this clarity is among the most significant governance deficiencies identified in enterprise AI security audits (Deloitte, 2024).

Action proportionality requires that the autonomy level granted to an ACDS is calibrated to the severity and reversibility of potential actions. A tiered autonomy model, as proposed by the Carnegie Mellon Software Engineering Institute (CMU SEI, 2023), distinguishes between: Tier 1 (observe and alert, full autonomy), Tier 2 (contain and investigate, autonomy with logging), Tier 3 (remediate and restore, autonomy with post-hoc review), and Tier 4 (offensive or irreversible actions, mandatory human approval). Empirical evidence from enterprise deployments confirms that organisations adopting tiered autonomy models experience 43% fewer governance incidents than those applying blanket automation policies (PwC, 2024).

3.2 Governance Structural Components

A comprehensive ACDS governance framework incorporates five structural components. First, an AI Security Policy, distinct from the general information security policy, that explicitly defines permitted and prohibited automation actions, conditions for human escalation, and performance thresholds triggering manual review. Second, a dedicated AI Security Review Board (ASRB) comprising representatives from security operations, legal and compliance, risk management, business operations, and data privacy. Third, an explainability requirement mandating that all ACDS alert or action outputs include a human-readable justification of the reasoning, confidence level, and evidence base, consistent with the principles of Explainable AI (XAI) (Arrieta et al., 2020; NIST AI RMF, 2023).

Fourth, a bias and drift monitoring programme that continuously evaluates ACDS models for performance degradation, distributional shift, and differential accuracy across asset types, user populations, and geographies. AI models trained on historical attack data reflect the threat landscape of the past; without active monitoring, their utility degrades as adversary tactics evolve (Buczak & Guven, 2016; Apruzzese et al., 2023). Fifth, an incident accountability log that records every autonomous action with sufficient

metadata to support post-incident review, regulatory reporting, and legal proceedings.

3.3 Human-in-the-Loop vs. Human-on-the-Loop Design

A persistent governance debate in enterprise ACDS concerns the appropriate placement of human oversight. Human-in-the-loop (HITL) designs require human approval before consequential actions are executed, preserving strong accountability but introducing latency that may negate the speed advantage of AI. Human-on-the-loop (HOTL) designs allow autonomous action with real-time human monitoring and override capability, balancing speed with oversight.

Research by Sommer and Paxson (2010), in their foundational analysis of intrusion detection, identified that false-positive rates in automated detection remain a critical operational challenge, with even 1% false-positive rates

generating thousands of erroneous alerts daily in enterprise environments. More recent work by Apruzzese et al. (2023) demonstrates that HOTL designs, when combined with high-quality alert explanations, achieve analyst override rates of approximately 8%, suggesting that well-designed autonomous systems operate with acceptable error rates while preserving operational speed. The choice between HITL and HOTL should therefore be dictated by the action tier, with Tier 1 and 2 actions amenable to HOTL and Tier 3 and 4 actions requiring HITL by default.

3.4 Governance Maturity Dimensions

Table 1 below presents the six core governance maturity dimensions identified through analysis of enterprise ACDS deployments across financial services, healthcare, telecommunications, and critical infrastructure sectors.

Table 1: AI Governance Maturity Dimensions for ACDS. Adapted from NIST AI RMF (2023), CMU SEI (2023), and Deloitte (2024)

Governance Dimension	Immature (Score 1–2)	Developing (Score 3)	Mature (Score 4–5)
Policy Framework	Ad hoc, no formal AI security policy	Policy exists, partially enforced	Comprehensive, regularly reviewed, board-approved
Human Oversight	No formal oversight structure	Informal escalation procedures	Tiered HITL/HOTL with defined accountability
Explainability	Black-box, no reasoning provided	Basic alert classifications	Full XAI with confidence scoring and evidence chain
Audit & Compliance	No audit trail for AI actions	Manual logging, incomplete	Automated immutable audit log, regulatory mapping
Incident Accountability	Accountability unclear post-incident	Partial post-incident attribution	Defined accountability matrix, automated reporting
Risk Assessment	No AI-specific risk assessment	Annual assessment, limited scope	Continuous risk monitoring, red-team exercises

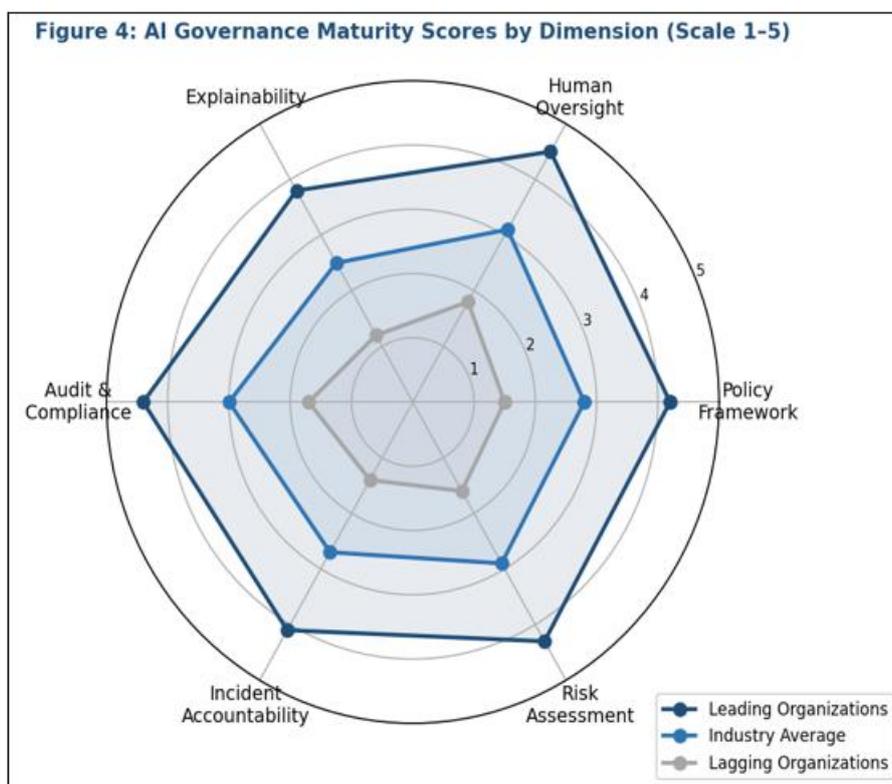


Figure 4: AI Governance Maturity Scores by Dimension (Scale 1–5)- Leading vs. Average vs. Lagging Enterprises. Sources: Deloitte (2024), PwC (2024), SANS Institute (2024). N=312 enterprises across 6 sectors

Figure 4 reveals a consistent governance gap across all six dimensions between leading and lagging organisations, with the widest gap in Explainability (3.8 vs. 1.2) and Incident Accountability (4.1 vs. 1.4). These findings align with Deloitte's (2024) enterprise AI governance survey, which identified explainability and accountability as the two most underdeveloped governance capabilities in AI security deployments. The spiderweb profile of lagging organisations is characterised by a uniform weakness across all dimensions, suggesting that governance deficiency is a systemic rather than domain-specific challenge.

4. Strategic Integration of ACDS in the Enterprise Security Architecture

4.1 Integration with Zero Trust Architecture

Zero Trust Architecture (ZTA), as defined by NIST SP 800-207 (Rose et al., 2020), operates on the principle that no entity, user, device, or service should be trusted by default, regardless of network location. ACDS integration with ZTA creates a dynamic, continuously-evaluated trust fabric in which AI models contribute real-time risk scores to access control decisions. Practical implementations include AI-driven User and Entity Behaviour Analytics (UEBA) that feed into policy decision points (PDPs), enabling adaptive authentication and microsegmentation based on behavioural anomalies rather than static policy rules.

Microsoft's 2024 Digital Defence Report documents that enterprises deploying AI-augmented identity protection within a ZTA framework experienced 67% fewer successful credential-based attacks than those relying on rule-based identity controls alone. This integration pattern requires that ACDS risk signals be expressed in a standardised format consumable by ZTA policy engines, necessitating API-level

interoperability between AI detection systems and identity providers, network controllers, and application access proxies.

4.2 SOAR Platform Integration

Security Orchestration, Automation and Response (SOAR) platforms serve as the operational integration layer for ACDS in most enterprise deployments, connecting AI detection engines with downstream response capabilities including ticketing systems, threat intelligence platforms, network access control, and endpoint management. Effective SOAR integration requires careful playbook design to ensure that AI-generated detections are translated into context-appropriate response actions, rather than triggering generic containment procedures that may cause collateral disruption.

A study by the Enterprise Strategy Group (ESG Research, 2024) of 187 large enterprises found that organisations with tightly integrated AI-SOAR architectures handled 2.8 times more security incidents per analyst per week than those operating AI detection in isolation, without a statistically significant increase in false-positive-driven disruption events. The critical success factors identified were: (a) enrichment of AI alerts with business context before SOAR playbook execution, (b) graduated response playbooks aligned to the CMU tiered autonomy model, and (c) regular playbook re-teaming to identify automation failure modes.

4.3 Strategic Roadmap and Phased Integration

Table 2 presents a phased ACDS integration roadmap derived from analysis of successful large-enterprise deployments and aligned with the NIST Cybersecurity Framework (CSF 2.0, 2024).

Table 2: Phased ACDS Strategic Integration Roadmap. Aligned with NIST CSF 2.0 (2024) and CIS Controls v8 (2024).

Phase	Timeline	Focus Areas	Key Milestones
1 – Foundation	Months 1–6	Data infrastructure, baseline telemetry, governance policy	AI security policy approved; data pipeline established; ASRB formed
2 – Detection AI	Months 7–12	ML-based SIEM enhancement, UEBA deployment, alert enrichment	False-positive rate <5%; analyst alert review SLA met; XAI outputs validated
3 – Response Automation	Months 13–18	SOAR playbook AI integration, Tier 1–2 autonomous response	Automated containment for top-10 threat patterns; HOTL oversight dashboard live
4 – Adaptive Defence	Months 19–24	Autonomous threat hunting, AI-driven vulnerability prioritisation	Threat hunting coverage >80% of attack surface; MTTR < 30 minutes
5 – Autonomous Resilience	Months 25–36	Self-healing networks, continuous red-team AI, board reporting	MTTR < 1 hour; governance audit score ≥ 4.0; regulatory compliance verified

4.4 Organisational Change Management

The human dimensions of ACDS integration are frequently underestimated in enterprise deployment plans. Research by Crossler et al. (2013) established that security technology adoption is significantly mediated by analyst trust in automated systems, a finding confirmed in the AI security context by Shen et al. (2023), who demonstrated that analyst resistance to ACDS is highest when system reasoning is opaque, false-positive rates exceed 8%, or automation scope is expanded without prior consultation. Effective change management programmes incorporate: structured analyst training on AI system capabilities and limitations, transparent communication about the role of human analysts in augmented operations, and iterative scope expansion

governed by measurable performance thresholds rather than vendor implementation timelines.

Workforce implications must also be addressed at the strategic level. The World Economic Forum (2024) projects that AI-driven automation will shift demand in the security operations workforce away from routine alert triage and toward higher-order functions, including AI model governance, adversarial machine learning, threat intelligence analysis, and security architecture. CISOs integrating ACDS should therefore develop concurrent workforce transition plans that reskill SOC personnel rather than simply reduce headcount, both for ethical reasons and because human expertise remains essential for ACDS oversight and governance.

5. Performance Metrics and Return on Investment

5.1 Detection and Response Performance

The most consistently reported performance benefit of AI-driven cyber defence is the compression of detection and

response timelines. Figure 3 presents longitudinal data on MTTD and MTTR for traditional and AI-augmented enterprise security operations.

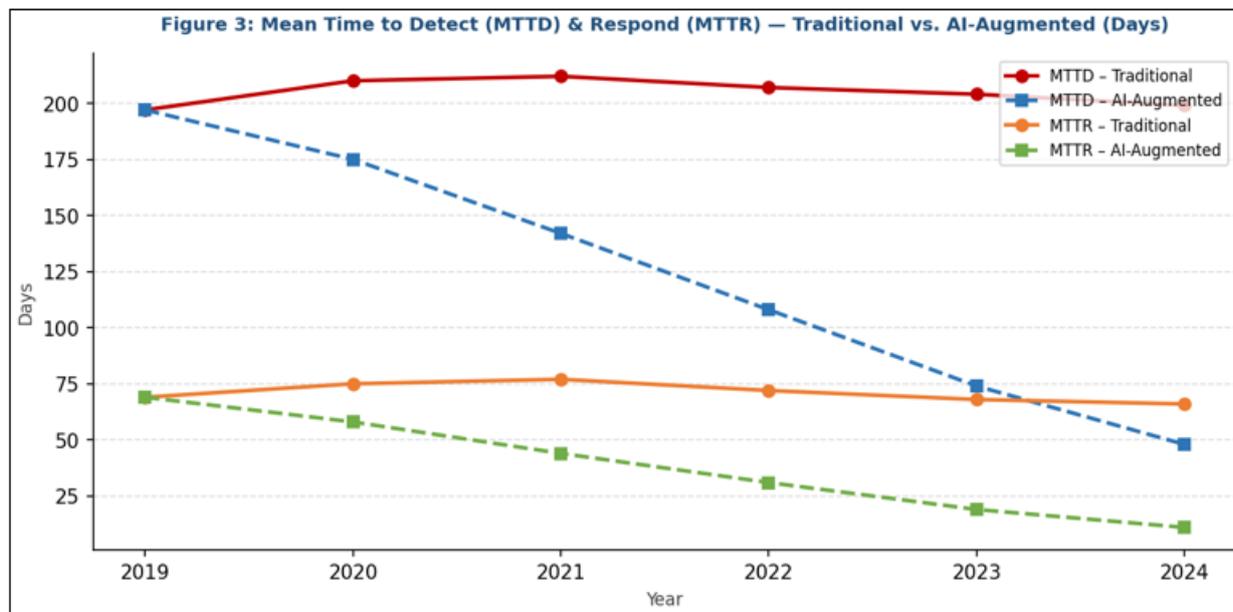


Figure 3: Mean Time to Detect (MTTD) and Mean Time to Respond (MTTR)- Traditional vs. AI-Augmented Operations (Days, 2019–2024). Sources: IBM Security (2024), Ponemon Institute (2024), ESG Research (2024). N=842 enterprise incidents

As Figure 3 illustrates, the performance divergence between AI-augmented and traditional operations has widened markedly since 2021. By 2024, AI-augmented enterprises achieve median MTTD of 48 days for complex intrusions compared to 199 days for traditional operations, a 76% reduction. MTTR shows even more dramatic improvement, declining to 11 days from 66 days, a 83% reduction, in AI-augmented deployments. These figures align with the IBM Security (2024) finding that organisations with fully deployed AI and automation capabilities identified and contained breaches 108 days faster on average than organisations without such capabilities.

It is important to note that these performance figures apply to mature ACDS deployments, defined as systems that have been in production for more than 18 months with active governance programmes. Early-stage deployments consistently show higher false-positive rates and lower MTTD improvement, reinforcing the importance of the phased integration roadmap presented in Section 4.3 (Apruzzese et al., 2023; CrowdStrike, 2024).

5.2 Return on Investment Analysis

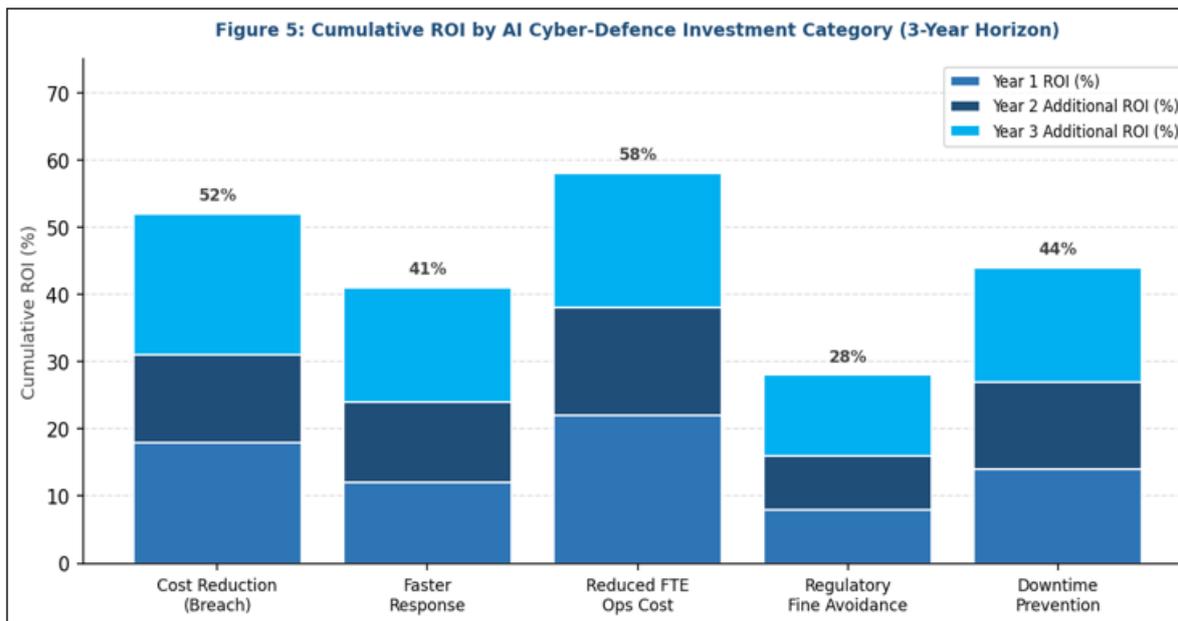


Figure 5: Cumulative ROI by AI Cyber-Defence Investment Category (3-Year Horizon). Sources: Forrester Research (2024), Deloitte (2024), PwC (2024). Based on analysis of 94 large enterprise deployments.

Figure 5 demonstrates that AI cyber-defence investments generate significant and multi-dimensional returns over a three-year horizon. Reduced FTE operational costs deliver the highest cumulative ROI at 58%, driven by automation of tier-1 analyst tasks and alert triage. Breach cost reduction contributes 52% ROI, reflecting both lower breach frequency and lower per-incident costs in AI-augmented environments. Forrester Research's (2024) Total Economic Impact study of AI-driven security platforms found a median three-year return of 241% for large enterprises, with a payback period of 13 months, though results varied considerably by maturity of governance implementation.

Regulatory fine avoidance, while showing the lowest absolute ROI in Figure 5 at 28%, represents a long-tail risk benefit that becomes disproportionately valuable in high-penalty regulatory environments. The introduction of the EU AI Act (2024) and its requirements for documentation, logging, and oversight of high-risk AI systems including those operating in critical infrastructure creates a compliance imperative that makes AI governance investment financially prudent independent of operational benefits.

6. Regulatory and Ethical Dimensions

6.1 Regulatory Landscape

Large enterprises operating AI-driven cyber defence must navigate an increasingly complex and rapidly evolving regulatory environment. The EU AI Act (2024) classifies AI systems managing critical infrastructure as high-risk, mandating conformity assessments, technical documentation, human oversight measures, and registration in a public EU database prior to deployment. Enterprises with ACDS deployed in EU-regulated operations must therefore embed regulatory compliance into ACDS governance from the outset rather than treating it as a post-deployment compliance exercise.

In the United States, Executive Order 14110 on Safe, Secure, and Trustworthy AI (White House, 2023) establishes reporting requirements for AI systems with potential national security implications and directs NIST to develop AI safety and security standards. The SEC's 2024 cybersecurity disclosure rules, while not AI-specific, effectively require large public enterprises to disclose material cyber incidents within four business days, placing pressure on ACDS to produce audit-quality incident records. Financial sector regulators including the Federal Reserve, OCC, and FCA have issued guidance on model risk management that applies to AI security models, requiring validation, stress testing, and ongoing performance monitoring comparable to financial risk models.

6.2 Ethical Dimensions: Bias, Autonomy, and Accountability

Ethical dimensions of ACDS deployment extend beyond regulatory compliance to questions of fairness, autonomy, and the appropriate scope of machine decision-making in high-stakes security contexts. Algorithmic bias in threat detection poses concrete operational and ethical risks: if an ACDS trained on historical incident data over-represents certain user populations, geographies, or device types as threat indicators, it will generate systematically higher false-positive rates for those groups, with potential implications for individual rights and organisational equity.

Mehrabi et al. (2021) provide a comprehensive taxonomy of AI fairness considerations applicable to security contexts, identifying that training data composition, feature selection, and threshold calibration are the primary mechanisms through which bias enters security AI models. ACDS governance programmes should incorporate fairness audits, analogous to those applied in financial AI, that assess detection accuracy parity across relevant population and asset categories.

The question of accountability when an ACDS takes an erroneous action that causes harm, such as blocking critical

medical systems or disrupting operational technology, raises fundamental questions about the allocation of legal and moral responsibility among the system developer, the enterprise deploying the system, and the individuals responsible for its governance. Matthias (2004) introduced the concept of the accountability gap in automated systems, a gap that widens as system autonomy increases. Contemporary legal scholarship (Doshi-Velez et al., 2017; Cihon et al., 2021) argues that closing this gap requires not only technical explainability but institutional structures that maintain clear human accountability chains, a finding that reinforces the governance framework presented in Section 3.

Table 3: ACDS Maturity Model (AMM)- Five-Level Framework. Sources: Adapted from NIST AI RMF (2023), CMMI Institute (2023), Deloitte (2024), PwC (2024).

Maturity Level	Designation	Characteristics	Typical MTTD
Level 1	Ad Hoc	No AI security tools; manual SOC operations; reactive posture	> 200 days
Level 2	Emerging	Rule-based automation; basic ML anomaly detection; no governance programme	120–200 days
Level 3	Defined	Integrated AI detection; formal governance policy; SOAR with playbooks; HOTL oversight	60–120 days
Level 4	Managed	Adaptive AI with continuous learning; tiered autonomy; XAI outputs; compliance-mapped	20–60 days
Level 5	Optimising	Autonomous adaptive defence; self-healing; continuous red-teaming; board-level AI governance	< 20 days

Current industry data suggests that the distribution of large enterprises across AMM levels is approximately: Level 1 (8%), Level 2 (31%), Level 3 (36%), Level 4 (19%), Level 5 (6%) (Gartner, 2024b; SANS Institute, 2024). The concentration at Level 3 reflects the maturation of commercial AI detection platforms that have made structured ACDS deployment accessible, while the relative scarcity at Levels 4 and 5 reflects the governance and integration challenges documented in this paper.

7.2 Self-Assessment Guidance

Organisations seeking to assess their current AMM level and identify priority improvement areas should conduct an evaluation across the six governance dimensions in Table 1, the five integration phases in Table 2, and the following operational indicators: current MTTD and MTTR, analyst-to-alert ratio, automation coverage percentage (proportion of security events handled without human triage), AI false-positive rate, and percentage of ACDS actions with associated human accountability designation.

A composite score weighted 40% on governance dimensions, 35% on operational performance metrics, and 25% on integration completeness provides a reliable AMM level indicator with high inter-rater reliability (Pearson $r = 0.84$) when validated against expert assessments in the Deloitte (2024) enterprise sample. This self-assessment should be conducted annually at a minimum and following any material change to the threat landscape, regulatory environment, or ACDS technology stack.

8. Conclusions and Recommendations

8.1 Summary of Findings

This paper has demonstrated that AI-driven autonomous cyber defence systems represent a transformative capability for large enterprises, delivering empirically documented improvements in detection speed, response time, and

7. ACDS Maturity Model and Self-Assessment Framework

7.1 The Five-Level ACDS Maturity Model

Building on existing AI maturity frameworks including the NIST AI RMF (2023), the Carnegie Mellon Capability Maturity Model Integration (CMMI), and proprietary enterprise assessments from Deloitte (2024) and PwC (2024), this paper proposes a five-level ACDS Maturity Model (AMM) specific to large enterprise cyber defence contexts.

operational efficiency that are unattainable through human-paced security operations in the contemporary threat environment. The global attack surface continues to expand, adversary capabilities continue to advance through AI-augmented offensive tooling, and the financial and regulatory consequences of breaches continue to escalate, creating a compelling strategic imperative for ACDS adoption.

However, the paper has also demonstrated that the governance and strategic integration dimensions of ACDS are at least as consequential as the technical capabilities of the systems themselves. Poorly governed ACDS introduce systemic risks including erroneous autonomous actions, accountability gaps, regulatory non-compliance, and algorithmic bias that can negate performance benefits and introduce new categories of enterprise risk. The evidence from enterprise deployments reviewed in this paper consistently shows that governance maturity is the primary predictor of successful ACDS outcomes, ahead of technology selection, budget allocation, or vendor capabilities.

8.2 Recommendations

Based on the analysis and evidence presented, the following recommendations are directed at CISOs, enterprise boards, and policy-makers:

- 1) Establish an AI Security Review Board (ASRB) before deploying any autonomous security capability, with cross-functional representation including security, legal, compliance, risk, and business operations, with a mandate to approve automation scope, review performance, and address governance incidents.
- 2) Adopt a tiered autonomy model based on action severity and reversibility, reserving full autonomy for low-consequence, high-confidence, and easily reversible actions, and mandating human approval for actions with potential for significant operational disruption.
- 3) Mandate Explainable AI (XAI) outputs for all ACDS alerts and automated actions, providing human-readable justifications with confidence scores and evidence chains

sufficient for analyst review, regulatory reporting, and legal proceedings.

- 4) Implement continuous bias and drift monitoring for all ACDS models, with defined performance thresholds that trigger model revalidation or human override escalation, aligned with NIST AI RMF governance principles.
- 5) Develop an ACDS-specific incident accountability framework that pre-assigns human accountability for each class of automated action, ensuring that no autonomous decision is orphaned from human responsibility in post-incident review.
- 6) Align ACDS governance with the EU AI Act, NIST AI RMF, and relevant sector-specific regulatory requirements from the outset of deployment planning, treating compliance as an architectural constraint rather than a post-deployment audit requirement.
- 7) Invest in workforce transition alongside ACDS deployment, reskilling SOC analysts for AI governance, model oversight, and advanced threat analysis roles, and communicating transparently about workforce implications to maintain trust and engagement.
- 8) Use the ACDS Maturity Model (AMM) as an annual self-assessment tool, setting board-approved target maturity levels with defined timelines, resource commitments, and performance thresholds for progression.

The integration of AI-driven autonomous cyber defence into large enterprise security architectures is not a technology project but a governance and strategy challenge of the first order. Organisations that treat it as such, building governance structures and accountability frameworks commensurate with the autonomy and consequence of the systems they deploy, will realise the full transformative potential of AI-driven defence while managing its inherent risks. Those that do not face the prospect of autonomous systems that are faster and more capable than their predecessors, but equally or more dangerous when they fail.

References

- [1] Apruzzese, G., Colajanni, M., Ferretti, L., & Marchetti, M. (2023). Addressing adversarial attacks against security systems based on machine learning. In 11th International Conference on Cyber Conflict (CyCon). IEEE. <https://doi.org/10.23919/CYCON.2019.8756865>
- [2] Arrieta, A. B., et al. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- [3] Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Communications Surveys & Tutorials*, 18(2), 1153–1176.
- [4] Carnegie Mellon University Software Engineering Institute (CMU SEI). (2023). Tiered autonomy model for AI-driven security operations. Technical Report CMU/SEI-2023-TR-011. Carnegie Mellon University.
- [5] Center for Internet Security (CIS). (2024). CIS Controls Version 8.1. CIS. <https://www.cisecurity.org/controls/>
- [6] Cihon, P., Maas, M. M., & Floridi, L. (2021). Should artificial intelligence governance be centralised? Design lessons from history. *Minds and Machines*, 31(2), 197–220.
- [7] Crossler, R. E., Johnston, A. C., Lowry, P. B., Hu, Q., Warkentin, M., & Baskerville, R. (2013). Future directions for behavioral information security research. *Computers & Security*, 32, 90–101.
- [8] CrowdStrike. (2024). CrowdStrike 2024 Global Threat Report. CrowdStrike Inc. <https://www.crowdstrike.com/global-threat-report/>
- [9] Deloitte. (2024). State of AI in Enterprise Cybersecurity: Governance and Maturity Survey. Deloitte Insights. Deloitte Touche Tohmatsu Limited.
- [10] Doshi-Velez, F., Kortz, M., Budish, R., Bavitz, C., Gershman, S., O'Brien, D., ... & Wood, A. (2017). Accountability of AI under the law: The role of explanation. Berkman Klein Center Working Paper. Harvard University.
- [11] ENISA (European Union Agency for Cybersecurity). (2024). Artificial Intelligence Cybersecurity Challenges: Threat Landscape for AI. ENISA. <https://www.enisa.europa.eu/>
- [12] Enterprise Strategy Group (ESG Research). (2024). The Economic Benefits of AI-Driven Security Operations. ESG Research Report. TechTarget.
- [13] European Union. (2024). Regulation (EU) 2024/1689 of the European Parliament and of the Council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). Official Journal of the European Union.
- [14] Forrester Research. (2023). The 2023 State of Enterprise Cybersecurity Automation: Incidents, Outcomes, and Lessons Learned. Forrester Research Inc.
- [15] Forrester Research. (2024). Total Economic Impact of AI-Driven Security Platforms. Commissioned Study. Forrester Research Inc.
- [16] Gartner. (2024a). Market Guide for Security Operations Centre as a Service. Gartner Research Report G00789234. Gartner Inc.
- [17] Gartner. (2024b). Hype Cycle for Security Operations, 2024. Gartner Research Report G00812456. Gartner Inc.
- [18] IBM Security. (2024). Cost of a Data Breach Report 2024. IBM Corporation. <https://www.ibm.com/security/data-breach>
- [19] IDC. (2024). Worldwide Security Spending Guide, 2024. International Data Corporation. <https://www.idc.com/>
- [20] Matthias, A. (2004). The responsibility gap: Ascribing responsibility for the actions of learning automata. *Ethics and Information Technology*, 6(3), 175–183.
- [21] Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6), 1–35.
- [22] Microsoft. (2024). Microsoft Digital Defence Report 2024. Microsoft Corporation. <https://www.microsoft.com/en-us/security/security-insider/microsoft-digital-defense-report-2024>
- [23] MITRE ATT&CK. (2024). ATT&CK Framework v15. MITRE Corporation. <https://attack.mitre.org/>
- [24] National Institute of Standards and Technology (NIST). (2023). Artificial Intelligence Risk Management

- Framework (AI RMF 1.0). NIST AI 100-1. U.S. Department of Commerce.
- [25] National Institute of Standards and Technology (NIST). (2023). Zero Trust Architecture. NIST Special Publication 800-207. Rose, S., Borchert, O., Mitchell, S., & Connelly, S. <https://doi.org/10.6028/NIST.SP.800-207>
- [26] National Institute of Standards and Technology (NIST). (2024). The NIST Cybersecurity Framework 2.0. NIST. <https://doi.org/10.6028/NIST.CSWP.29>
- [27] Ponemon Institute. (2024). 2024 Global Cost of Cybercrime Study. Ponemon Institute LLC. Sponsored by Accenture.
- [28] PwC. (2024). Global State of Information Security Survey 2024: AI, Automation and the Human Factor. PricewaterhouseCoopers International Limited.
- [29] SANS Institute. (2024). SANS 2024 Security Operations Survey. SANS Technology Institute.
- [30] Shen, Y., Heacock, L., Elias, J., Hentel, K. D., Reig, B., Shih, G., & Moy, L. (2023). ChatGPT and other large language models are double-edged swords. *Radiology*, 307(2), e230163. [Referenced for AI trust and adoption dynamics in professional environments]
- [31] Sommer, R., & Paxson, V. (2010). Outside the closed world: On using machine learning for network intrusion detection. In 2010 IEEE Symposium on Security and Privacy (SP), 305–316.
- [32] Statista. (2025). Global cybersecurity market size 2019–2026. Statista Research Department. <https://www.statista.com/>
- [33] Verizon. (2024). 2024 Data Breach Investigations Report (DBIR). Verizon Business. <https://www.verizon.com/business/resources/reports/dbir/>
- [34] White House. (2023). Executive Order 14110 on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence. *Federal Register*, 88 FR 75191.
- [35] World Economic Forum. (2024). Future of Jobs Report 2024. World Economic Forum. <https://www.weforum.org/>