



Hochschule Karlsruhe  
University of Applied Science

Fakultät für Informatik und Wirtschaftsinformatik  
Wirtschaftsinformatik

Masterthesis

Generative AI for Security Automation in Hyperscale Cloud Platforms

Von	VON
Matrikelnr.	0000
Arbeitsplatz	ORT
Erstbetreuer	PROF1
Zweitbetreuer	PROF2
Abgabetermin	DATUM

Karlsruhe, DATUM

Vorsitzender des Prüfungsausschusses



## Declaration of Authorship

I, Daniel VERA GILLIARD, in lieu of an oath that I have written the Master's thesis presented here independently and exclusively using the literature and other aids provided. The thesis has not been submitted in the same or a similar form to any other examination authority for the award of an academic degree.

Signed:

---

Date:

---



*“Thanks to my solid academic training, today I can write hundreds of words on virtually any topic without possessing a shred of information, which is how I got a good job in journalism.”*

Dave Barry



HOCHSCHULE KARLSRUHE

# *Abstract*

Faculty Name  
Business Information Systems

Master of Business Information Systems

**Generative AI for Security Automation in Hyperscale Cloud Platforms**

by Daniel VERA GILLIARD

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...





## *Acknowledgements*

The acknowledgments and the people to thank go here, don't forget to include your project advisor...



# Contents

<b>Declaration of Authorship</b>	<b>iii</b>
<b>Abstract</b>	<b>vii</b>
<b>Acknowledgements</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Instead of an introduction . . . . .	1
<b>2 Background and Related Work</b>	<b>3</b>
2.1 Foundational Concepts in Cloud Computing . . . . .	3
2.2 Foundational Concepts in Generative AI . . . . .	3
2.3 State of Cloud Provider Ecosystems . . . . .	3
2.4 Literature Review . . . . .	3
Methodology . . . . .	4
AI-Driven Security Approaches . . . . .	4
GenAI Security Scoping Matrix . . . . .	5
GenAI Security Frameworks . . . . .	6
Agent-Based Approaches . . . . .	10
Security Risks . . . . .	11
Balance of Automation and Human Oversight . . . . .	16
Summary Literature review . . . . .	17
2.5 Research Gaps . . . . .	17
<b>A Appendix Title Here</b>	<b>19</b>



# List of Figures



# List of Tables





# Listings



# List of Abbreviations

<b>ABAC</b>	<b>A</b> tttribute- <b>B</b> ased <b>A</b> ccess <b>C</b> ontrol
<b>ACSC</b>	<b>A</b> ustralian <b>C</b> yber <b>S</b> ecurity <b>C</b> entre
<b>AI</b>	<b>A</b> rtificial <b>I</b> ntelligence
<b>AI RMF</b>	<b>A</b> rtificial <b>I</b> ntelligence <b>R</b> isk <b>M</b> anagement <b>F</b> ramework
<b>APP</b>	<b>A</b> ustralian <b>P</b> rivacy <b>P</b> inciples
<b>CIA</b>	<b>C</b> onfidentiality, <b>I</b> ntegrity, and <b>A</b> vailability
<b>CSP</b>	<b>C</b> loud <b>S</b> ervice <b>P</b> rovider
<b>DTA</b>	<b>D</b> igital <b>T</b> ransformation <b>A</b> gency
<b>GenAI</b>	<b>G</b> enerative <b>A</b> rtificial <b>I</b> ntelligence
<b>LLM</b>	<b>L</b> arge <b>L</b> anguage <b>M</b> odel
<b>MLOps</b>	<b>M</b> achine <b>L</b> earning <b>O</b> perations
<b>MTTD</b>	<b>M</b> ean <b>T</b> ime to <b>D</b> etect
<b>MTTR</b>	<b>M</b> ean <b>T</b> ime to <b>R</b> esolve
<b>RAG</b>	<b>R</b> etrieval- <b>A</b> ugmented <b>G</b> eneration
<b>RPO</b>	<b>R</b> ecovery <b>P</b> oint <b>O</b> bjectives
<b>RTO</b>	<b>R</b> ecovery <b>T</b> ime <b>O</b> bjectives
<b>SOC</b>	<b>S</b> ecurity <b>O</b> perations <b>C</b> enter
<b>SRM</b>	<b>S</b> hared <b>R</b> esponsibility <b>M</b> odel
<b>WAF</b>	<b>W</b> eb <b>A</b> pplication <b>F</b> irewalls
<b>ZTA</b>	<b>Z</b> ero <b>T</b> rust <b>A</b> rchitecture



# Physical Constants

Speed of Light  $c_0 = 2.997\,924\,58 \times 10^8 \text{ m s}^{-1}$  (exact)



# List of Symbols

$a$	distance	m
$P$	power	W (J s <sup>-1</sup> )
$\omega$	angular frequency	rad





*For/Dedicated to/To my...*



## Chapter 1

# Introduction

### 1.1 Instead of an introduction

First of all: The introduction should be short!

State the problem, describe the organization and structure of the document and that's it. Anything more than 3 pages needs justification.



## Chapter 2

# Background and Related Work

### 2.1 Foundational Concepts in Cloud Computing

TBD

### 2.2 Foundational Concepts in Generative AI

TBD

### 2.3 State of Cloud Provider Ecosystems

TBD

### 2.4 Literature Review

The convergence of Generative Artificial Intelligence (GenAI) and hyperscale cloud platforms presents a rapidly evolving frontier for cybersecurity. As organizations increasingly rely on complex cloud environments, the scale and sophistication of threats necessitate advanced automation capabilities. GenAI offers promising avenues for enhancing security posture through intelligent automation, but its integration also introduces challenges and risks. A comprehensive understanding of the existing research landscape is therefore essential to identify established practices, emerging trends, and critical gaps in knowledge.

This literature review synthesizes current academic and industry contributions pertinent to the application of GenAI for security automation within hyperscale and multi-cloud contexts. It begins by outlining the methodology employed to select and analyze relevant works. Subsequently, the review dives into several key thematic areas: the evolution from traditional security methods to AI-driven approaches, specific frameworks for scoping and managing GenAI security implementations, architectural patterns and techniques for security automation, a critical examination of the unique security risks associated with GenAI itself, risk management frameworks and the ongoing crucial discussion regarding the necessary balance between automation and human oversight.

By examining these facets, this review aims to provide a foundational understanding of the state-of-the-art, highlighting both the potential of GenAI in cloud security automation and the significant considerations that must be addressed for its responsible and effective deployment. This synthesis will inform the subsequent research presented in this thesis.

## Methodology

This literature review followed a structured approach to identify relevant publications, focusing on peer-reviewed articles addressing GenAI applications in hyper-scale cloud security published primarily within the last five years. The search utilized academic databases with key search terms related to generative AI, cloud security automation, hyperscale platforms, and multi-cloud orchestration. Papers were selected based on their relevance to:

- GenAI applications specifically in cloud security contexts
- Hyperscale or multi-cloud environments
- Technical solutions for security automation
- Empirical evidence or theoretical frameworks with substantial methodological rigor

The selection process involved initial screening of titles and abstracts followed by full-text review of promising papers. The analysis employed a thematic approach, identifying recurring concepts, methodological approaches, and gaps in existing research. Particular attention was paid to identifying the theoretical foundations underpinning GenAI applications in security contexts, empirical evidence of effectiveness, and limitations of current approaches.

## AI-Driven Security Approaches

The landscape of cloud security is undergoing a significant transformation, shifting from traditional, often reactive methods towards more proactive and adaptive strategies powered by Artificial Intelligence (AI), particularly Generative AI (GenAI). This evolution marks a move beyond basic anomaly detection towards sophisticated security postures capable of learning from and responding dynamically to novel threat vectors in complex cloud environments.

Foundational work by Khanna [1] explores the integration of GenAI into cloud security, outlining its core applications. According to Khanna, modern GenAI implementations focus on key capabilities such as processing vast amounts of data (e.g., log entries, network packets) for advanced anomaly detection and threat intelligence, enabling automated response mechanisms that dynamically adjust security protocols, and facilitating predictive security measures to forecast potential vulnerabilities. While highlighting these advancements, Khanna also acknowledges inherent challenges, including the need for large datasets and mitigating potential adversarial manipulation [1].

Building upon these foundational capabilities, the integration of GenAI represents a significant leap beyond conventional rule-based security systems. Research indicates that GenAI enhances security automation, particularly within multi-cloud and hybrid architectures. It allows systems to adapt infrastructure dynamically in response to varying traffic patterns and implements AI-powered defenses against continuously evolving cyber threats. This adaptive capability directly addresses persistent challenges in cloud security related to optimizing workload distribution, ensuring performance, and managing costs effectively [2].

Furthermore, recent studies underscore the practical impact of integrating GenAI with established cloud-native security tools. Patel et al. [3] demonstrate how layering GenAI onto platforms like AWS GuardDuty and Google Cloud Security Command

Center significantly boosts automated threat detection, enables real-time incident response, and improves comprehensive vulnerability management across distributed cloud infrastructures. Their work provides empirical evidence, citing examples like Netflix and JPMorgan Chase, which reported measurable improvements in detection accuracy and notable reductions in security incidents following the adoption of GenAI-driven security automation strategies [3]. This convergence of GenAI with existing security frameworks highlights its potential to transform security operations centers (SOCs) by enhancing both efficiency and effectiveness.

### GenAI Security Scoping Matrix

With the introduction of the "Generative AI Security Scoping Matrix," by AWS, a structured and comprehensive framework for assessing security requirements based on the type of GenAI deployment is provided[4]. This framework aids organizations in evaluating their security posture, identifying potential vulnerabilities, and implementing appropriate controls throughout the AI lifecycle[4]. While core security disciplines remain essential, the matrix specifically helps address the unique risks and additional considerations introduced by generative AI workloads[4]. It classifies implementations into five scopes, representing increasing levels of ownership and control[5]:

1. **Consumer apps (Scope 1):** Utilising public third-party GenAI services (e.g., public chatbots), often with no cost or paid access, where the business does not own or see the training data or model and cannot modify it. Interaction occurs via APIs or direct application use according to provider terms[5][4]. This falls under the "Buying" category of GenAI usage[4].
2. **Enterprise apps (Scope 2):** Employing third-party enterprise applications (e.g., Amazon Q) with embedded GenAI features, involving an established business relationship with the vendor[5][4]. This also falls under the "Buying" category[4].
3. **Pre-trained models (Scope 3):** Building custom applications that integrate with existing third-party foundation models (e.g., Amazon Bedrock base models) via APIs[5][4]. This represents the start of the "Building" category[4].
4. **Fine-tuned models (Scope 4):** Refining an existing third-party foundation model by fine-tuning it with business-specific data, resulting in a new, specialized model (e.g., Amazon Bedrock customized models)[5][4]. This is part of the "Building" category[4].
5. **Self-trained models (Scope 5):** Constructing and training a GenAI model from the ground up using proprietary or acquired data, implying full ownership of the model (e.g., using Amazon SageMaker)[5][4]. This represents the highest level of ownership in the "Building" category[4].

This matrix serves as a mental model[4] that helps security teams prioritize focus areas by identifying five key security disciplines whose requirements vary across the deployment scopes: governance and compliance, legal and privacy, risk management, controls, and resilience[5][4]. As organizations move across the scopes from consuming consumer apps (Scope 1) towards building self-trained models (Scope 5), the demands within these disciplines evolve significantly. For instance, governance requirements escalate from basic compliance with terms of service to comprehensive frameworks for model development and monitoring. Similarly, risk management priorities shift, potentially focusing on prompt injection for pre-trained models (Scope

3) versus data poisoning or model extraction for fine-tuned or self-trained models (Scopes 4 and 5). The necessary security controls also transition from emphasizing access policies in lower scopes to implementing technical safeguards like adversarial testing and output filtering in higher scopes. Resilience planning also adapts based on the application's criticality. By mapping their GenAI activities to this matrix, organizations can systematically assess risks and apply appropriate security measures tailored to their specific implementation context.

### GenAI Security Frameworks

A notable contribution to the field is the SecGenAI framework, which provides a comprehensive approach to securing cloud-based Generative AI (GenAI) applications, with particular attention to Retrieval-Augmented Generation (RAG) systems within the context of Australian Critical Technologies of National Interest[6]. This framework addresses the unique security challenges introduced by the rapid advancement of GenAI technologies[6].

SecGenAI is structured around three core pillars: functional, infrastructure, and governance requirements[6]. It integrates an end-to-end security analysis to generate detailed specifications. These specifications emphasize critical areas such as data privacy, secure deployment methodologies, and the implementation of shared responsibility models between cloud service providers and users[6]. The framework's development addresses key questions surrounding GenAI security, the requirements for Confidentiality, Integrity, and Availability (CIA triad), specific RAG implementation options, constraints within the Australian regulatory landscape, and alignment with ethical principles[6].

A key aspect of SecGenAI is its alignment with established Australian guidelines, including the Australian Privacy Principles (APP) [7], particularly APP 11 concerning the security of personal information, Australia's AI Ethics Principles [7], and guidance from the Australian Cyber Security Centre (ACSC) and the Digital Transformation Agency (DTA)[6]. This alignment ensures that security measures incorporate regulatory compliance without hindering operational efficiency[6]. The framework specifically targets GenAI-specific threats that often evade traditional security countermeasures. These threats include data leakage (potentially revealing sensitive training or knowledge base data), various adversarial attacks (such as prompt injection, data poisoning, and model inversion techniques designed to extract underlying data or manipulate outputs), jailbreaking attacks (bypassing content restrictions), and the potential for GenAI to generate insecure code or be misused by threat actors[6].

SecGenAI proposes a multi-layered security strategy combining advanced machine learning techniques with robust security measures.

Functional Security Requirements:

- **Identity and Access Management:** Utilizing continuous and adaptive authentication mechanisms, alongside Attribute-Based Access Control (ABAC), to manage user access dynamically based on behavior and context[6].
- **Data Confidentiality and Integrity:** Employing techniques like homomorphic encryption to process encrypted data, data masking and tokenization to protect sensitive information, and data integrity verification using methods like hashing and artificial fingerprinting[6].



- **Model Security:** Implementing adversarial attack mitigation, encrypting model parameters, and ensuring secure model training protocols, potentially using differential privacy or federated learning[6].

#### Infrastructure Security Requirements:

- Implementing sandboxed environments (e.g., using containerization or virtualization) often within dedicated cloud availability zones and virtual private networks[6].
- Securing database connections using read replicas, strict IAM policies, and robust encryption methods (e.g., AWS KMS or CloudHSM)[6].
- Establishing stringent network security settings, segregating internal and external connections, and utilizing security groups[6].
- Deploying external attack prevention mechanisms like Web Application Firewalls (WAF) and DDoS mitigation services, potentially integrated with data processing and monitoring tools (e.g., AWS Kinesis, Glue, Athena, Grafana)[6].
- Ensuring robust data backup and disaster recovery strategies, considering Recovery Time Objectives (RTO) and Recovery Point Objectives (RPO)[6].

#### Governance Requirements:

- Adhering to AI governance principles (based on the ISO 38500 Evaluate-Direct-Monitor cycle [8]) covering fairness, accountability, content traceability (e.g., watermarking), data protection, regular audits, reliability, user consent, third-party risk management, transparency, and legal compliance[6].
- Clearly defining roles and responsibilities through an AI-specific Shared Responsibility Model (SRM) for cloud environments, outlining obligations for Cloud Service Providers (CSPs), customers, and areas of shared duty[6].

By addressing these functional, infrastructure, and governance aspects in detail, the SecGenAI framework provides actionable strategies for the secure implementation and operation of cloud-based GenAI systems, fostering innovation while safeguarding national interests and enhancing the overall reliability and trustworthiness of these transformative technologies[6].

As organizations increasingly adopt multi-cloud strategies, effective policy orchestration across diverse environments becomes critical for maintaining consistent security postures and operational efficiency[9]. A comprehensive analysis of unified AI and cloud platforms published in 2024 examines architectural frameworks and integration patterns that enable the convergence of AI tools, machine learning operations (MLOps), data processing systems, and workflow orchestration within cloud-native environments, primarily focusing on transforming process automation and decision systems[9]. This research investigates how these unified platforms address challenges like managing distributed AI workloads, ensuring real-time processing, and maintaining regulatory compliance[9].

This research identifies three key innovations within these unified platforms with significant implications for advanced automation, including security automation:

- **Federated AI implementations:** These allow organizations to train AI models across distributed nodes or cloud boundaries while preserving data sovereignty and privacy, as sensitive data does not need to be centralized. This approach also reduces network bandwidth requirements and utilizes techniques like secure aggregation protocols and differential privacy[9].

- **Real-time data processing architectures:** Leveraging advanced streaming technologies (like Apache Kafka or Flink), in-memory processing, and real-time analytics engines, these architectures enable sub-second decision-making and immediate response to events (such as potential threats) by efficiently handling continuous data streams with high reliability and fault tolerance[9].
- **Multi-cloud integration patterns:** These patterns establish standardized interfaces, communication protocols, service discovery mechanisms, load balancing, and security controls to ensure seamless and consistent operation, including policy enforcement, across different cloud providers. Hybrid cloud deployment strategies are also highlighted, intelligently distributing workloads between on-premise and cloud resources based on performance, cost, and compliance needs, managed by sophisticated orchestration[9].

These architectural approaches, integrating MLOps frameworks for lifecycle management, robust data processing, workflow orchestration, and advanced capabilities like federated learning and real-time analytics within a multi-cloud context, provide the necessary foundation upon which advanced solutions like GenAI-driven security automation can be built across heterogeneous cloud environments[9]. While the analyzed paper focuses broadly on process automation and decision systems, the detailed exploration of scalable, resilient, governable, and interoperable architectures directly supports the implementation of sophisticated, automated security measures in complex multi-cloud settings[9].

For organizations operating containerized workloads across multiple clusters, particularly in multi-domain architectures involving different administrative entities, research from 2023 proposes an automated approach for generating network security policies in Kubernetes deployments[10]. Manually configuring security in such environments is complex, often leading to inconsistencies between policies defined in different clusters and requiring domain administrators to possess knowledge about other domains' configurations (like service locations or IP addresses), which is not always feasible[10]. This approach addresses two critical challenges in multi-cluster security: reducing the configuration errors commonly made by human administrators and creating transparent cross-cluster communications without requiring extensive information sharing between domains[10].

The proposed solution involves a top-level entity named the "Multi-Cluster Orchestrator"[10]. This orchestrator acts as a central management point, receiving inputs from managers of different domains[10]. These inputs include:

- A description of each domain's structure (listing clusters and exposed services with their details)[10].
- High-level security requirements specifying allowed communications (e.g., between services within the same domain, services in different domains, or services and external IP addresses)[10]. These requirements can be defined using an extended YAML syntax with special labels ('service', 'cluster', 'domain') that abstract away low-level details[10].

Based on these inputs, the Multi-Cluster Orchestrator refines the high-level requirements into concrete configurations through a two-step process[10]:

1. It generates a "Global Configuration" that tracks communication pairs between services and required links between clusters, optimizing the overall cluster mesh setup[10].

2. It derives "Single Configurations" for each individual cluster, containing the specific parameters needed to connect the cluster to the mesh (e.g., using technologies like Cilium Cluster Mesh), the Kubernetes Network Policies to enforce the desired security rules, and commands to create local service entries that enable transparent name resolution for services located in external clusters[10].

The implementation, known as Multi-Cluster Orchestrator (developed in Java with REST APIs), demonstrates how automated policy generation can improve security consistency across distributed environments while reducing the cognitive load on security administrators by handling the complexity of multi-domain interactions transparently[10]. This research is particularly relevant for hyperscale cloud platforms and organizations that utilize container orchestration technologies like Kubernetes to manage numerous workloads across multiple clusters, potentially spanning different regions, availability zones, or administrative boundaries[10].

Another approach to security automation in the context of policies involves the use of digital twins for validating security policies before deployment in production environments[11]. This approach utilizes an emulation system specifically designed to create high-fidelity digital replicas of target IT infrastructures[11]. These digital twins replicate key functionalities of the corresponding physical or virtual systems, allowing security teams to play out complex security scenarios, such as intrusion attempts and defense responses, within a safe and controlled environment[11]. This capability avoids impacting operational workflows on the real-world infrastructure[11].

The digital twin approach, as detailed in the research by Hammar and Stadler, facilitates a closed-loop learning process for developing and refining security policies[11]. The process involves several key steps:

1. **Create Digital Twin:** A digital twin of the target infrastructure is generated using an emulation system built on virtualization technologies like Docker containers, virtual links, and virtual switches[11].
2. **Run Scenarios & Collect Data:** Security scenarios, involving emulated attackers, defenders, and client populations, are executed within the digital twin[11]. During these runs, detailed system measurements and logs are collected via monitoring agents that push metrics to data pipelines (e.g., using Kafka and Spark)[11].
3. **Model & Learn:** The collected data and statistics are used to instantiate simulations, often modeled as Markov decision processes [11]. Reinforcement learning techniques are then applied to these simulations to learn potentially optimal security policies[11].
4. **Validate & Iterate:** The performance of the learned policies is then rigorously evaluated back within the high-fidelity digital twin environment[11].

This methodology provides continuous, iterative feedback and improvement cycles, as the results from validation can inform further scenario runs and learning phases, enhancing policy effectiveness over time[11]. The authors demonstrate this by applying the approach to an intrusion response scenario, showing that the digital twin provided the necessary evaluative feedback to learn near-optimal policies that outperformed baseline systems like the SNORT IDPS[12]. This represents a significant advancement in validation mechanisms, particularly relevant for potentially complex GenAI-driven security automation strategies, by bridging the gap between simulation-based learning and real-world applicability[11].

Regarding policies, ensuring the trustworthiness and accuracy of GenAI-generated security policies and responses remains a significant challenge. The already mentioned SecGenAI framework demonstrates how advanced machine learning techniques can be combined with robust security measures to enhance the reliability of GenAI systems while maintaining compliance with regulatory requirements.[6] As described, this approach integrates continuous validation processes throughout the AI lifecycle, from model development to deployment and monitoring, creating multiple checkpoints that verify the integrity and effectiveness of security responses. By emphasizing explainability alongside accuracy, the framework addresses one of the primary concerns associated with GenAI applications in security contexts: the "black box" nature of complex models.[6]

While not specifically focused on cloud security, research on GenAI applications in the energy sector offers transferable insights into implementation approaches for complex operating environments. This comprehensive literature review identifies how GenAI enhances productivity through data creation, forecasting, optimization, and natural language understanding, while also addressing challenges such as hallucinations, data biases, privacy concerns, and system errors [13]. The proposed solutions—including improving training data quality, implementing system fine-tuning processes, establishing human oversight mechanisms, and deploying robust security measures—provide a valuable framework for GenAI implementations in cloud security contexts. These approaches are particularly relevant for hyperscale environments where scale and complexity amplify both the benefits and risks of GenAI adoption [13].

### Agent-Based Approaches

A recent paper from 2024 introduces and validates the concept of employing Generative AI (GenAI)-driven agentic workflows to achieve comprehensive security automation, particularly in complex modern environments. A notable example is the DevSecOps Sentinel system[14], specifically designed to address the mounting security challenges inherent in modern software supply chains. Challenges coming from microservices, containerization, and cloud-native architectures that often outpace traditional DevSecOps practices[14].

The DevSecOps Sentinel system exemplifies this approach by utilizing intelligent agents integrated into automated workflows. These agents are powered by advanced GenAI models, such as Large Language Models (LLMs) enhanced with Retrieval-Augmented Generation (RAG), enabling sophisticated analysis capabilities[14]. Key characteristics of these agents include:

- **Autonomy:** Operating independently based on predefined goals and policies.
- **Reactivity:** Responding in real-time to environmental changes like new vulnerability disclosures.
- **Proactivity:** Taking initiative, such as preemptively scanning for risks or suggesting improvements[14].

These agents execute critical security tasks throughout the software development lifecycle, including:

- **Automated Vulnerability and Impact Analysis:** Leveraging GenAI to analyze code, dependencies (tracked via SBOMs), and infrastructure configurations for potential threats, assessing their potential impact in context[14].

- **Adaptive Compliance and Release Gating:** Enforcing security policies and compliance requirements dynamically, acting as automated checks before deployment[14].
- **Predictive Security:** Utilizing AI to identify potential future risks based on historical data and emerging threat patterns[14].

The implementation and testing of DevSecOps Sentinel demonstrate several key points relevant to broader security automation:

1. **Viability for Complexity:** Agentic workflows powered by GenAI are shown to be a viable and effective method for tackling the intricate and rapidly evolving security issues found in modern, distributed systems[14].
2. **Synergy of AI and Agents:** The integration of GenAI's deep analysis capabilities with the autonomous, proactive nature of agentic systems offers a powerful paradigm for strengthening organizational security posture[14]. While Sentinel focuses on the supply chain, the principle applies broadly to automating security operations in complex cloud environments.
3. **Measurable Improvements:** Such systems can contribute to building and deploying software that is simultaneously faster, safer, and more reliable. The DevSecOps Sentinel study reported significant quantitative improvements in key security and operational metrics, including reduced Mean Time to Detect (MTTD) and Resolve (MTTR) for vulnerabilities, lower false positive rates, increased compliance pass rates, higher deployment frequency, and reduced change failure rates[14].

This approach, exemplified by DevSecOps Sentinel, highlights a promising direction for leveraging GenAI to automate and enhance security functions, moving beyond traditional limitations to offer more adaptive, context-aware, and efficient security management in demanding environments like hyperscale clouds.

## Security Risks

The increasing integration of Generative Artificial Intelligence (GenAI) into various domains, including cybersecurity, presents significant opportunities but also introduces complex and multifaceted risks. Insights from recent literature reviews highlight these emerging challenges. A systematic literature review by Nyoto et al. [15], analyzing 17 relevant studies according to PRISMA 2020 guidelines **page\_prisma\_2021**, identifies several significant cybersecurity threats stemming primarily from the irresponsible application of GenAI technology. Complementing this, Surathunmanun et al. [13], while reviewing GenAI in the energy sector, outline key challenges that possess direct and critical relevance to security implementations, particularly within cloud environments reliant on third-party models and data. Synthesizing these findings provides a comprehensive overview of the risks:

- **Enhanced Malicious Content Generation and Misuse:** GenAI significantly lowers the barrier for creating sophisticated malicious content and tools. It can be abused to generate highly personalized and convincing phishing messages and social engineering tactics, increasing their effectiveness even with minimal target information [15]. Furthermore, GenAI facilitates the creation of effective ransomware and diverse forms of malware, potentially empowering individuals



with limited coding expertise to launch attacks [15]. Beyond typical malware, it can also generate other executable attack code, such as SQL injection scripts [15]. This potential for misuse is a major concern, where uncontrolled access or improper application can lead to significant harm [13]. This includes leveraging GenAI to bypass security controls through techniques like prompt injection or jailbreaking, as highlighted in literature concerning Large Language Models. [13].

- **Information Integrity:** GenAI poses substantial risks to information integrity. It enables the creation of highly realistic deepfake audio and video content, often without clear legal frameworks or consent, leading to potential fraud, manipulation, reputational damage, and the spread of disinformation [15]. Concurrently, GenAI models are prone to generating plausible but factually incorrect or nonsensical information, known as hallucinations [13], [15]. This issue often arises from poor quality training data or suboptimal parameter settings [13] and can be exacerbated by data poisoning during the model training phase [15]. In a security context, hallucinations can manifest as faulty threat analyses, incorrect vulnerability assessments, or misleading security recommendations [13]. Compounding this is the issue of data bias, where biases inherent in training data or introduced during feature selection lead to skewed or unfair outputs [13]. For security applications, this could result in certain threat types being consistently overlooked or specific user groups being unfairly flagged, thereby undermining the reliability of automated systems [13]. These challenges are often exacerbated by the inherent 'black box' nature of many LLMs, characterized by their complexity, lack of transparency in internal decision-making, and limited explainability, making it difficult to fully diagnose or prevent issues like hallucinations or bias [16].
- **Data Privacy, Security Vulnerabilities, and Intellectual Property:** The foundation of GenAI models—vast datasets—introduces significant privacy and security risks. Models are often trained on data scraped without explicit consent, potentially including sensitive personal information or copyrighted material [15]. User interactions and prompts can also be incorporated into training data, leading to potential data leakage and privacy violations [15]. This raises substantial intellectual property concerns and challenges compliance with regulations like GDPR [15]. The lack of transparency and control over how data is utilized presents considerable privacy risks [15]. Furthermore, insecure data handling practices can create security vulnerabilities [13]. Specific risks associated with LLMs, often used in cloud-hosted GenAI services, include inference attacks, data extraction attacks, data poisoning, supply chain vulnerabilities [13] and vulnerabilities to adversarial attacks stemming from the models complex and often opaque nature [16].
- **Systemic and Operational Risks:** Beyond content generation and data issues, GenAI systems can introduce operational risks. Logical inconsistencies within the model or unforeseen external events can cause GenAI systems to produce errors or fail entirely [13]. In automated security workflows operating in cloud environments (e.g., incident response, configuration management), such errors could propagate rapidly, leading to service disruptions, critical misconfigurations, or a failure to respond effectively to genuine threats [13].

These diverse risks, spanning malicious misuse, information integrity compromises, privacy violations, intellectual property infringements, and operational failures, underscore the critical need for robust countermeasures and responsible governance. Addressing these challenges necessitates comprehensive approaches, including rigorous data governance frameworks, cross-verification of GenAI outputs, continuous model monitoring and updating, incorporating human-in-the-loop validation processes, implementing strong security measures [13] and architectures like Zero Trust [13], [16], establishing clear ethical guidelines, and potentially developing new regulations specific to GenAI development and deployment [15]. Ensuring the responsible use of GenAI is paramount to harnessing its benefits while mitigating the significant emerging cybersecurity challenges, particularly in sensitive contexts like cloud security where the consequences of unreliable or misused AI can severely impact organizational risk posture and operational integrity [13].

Another significant challenge in implementing GenAI for security automation is the comprehensive identification and management of the unique risks these systems introduce, which differ significantly from traditional software risks. The NIST Artificial Intelligence Risk Management Framework (AI RMF 1.0) [17] provides a structured, voluntary approach to address these challenges.

The AI RMF defines an AI system as an "engineered or machine-based system that can, for a given set of objectives, generate outputs such as predictions, recommendations, or decisions influencing real or virtual environments" [17, p.1]. It acknowledges that while AI offers transformative potential, it also poses distinct risks due to factors like data dependency, complexity, opacity, and the socio-technical context of deployment [17].

In the paper, the NIST describes some key points relevant to GenAI Security Risks in Cloud Computing.

1. **Unique AI Risk Landscape:** The framework highlights that AI risks differ from traditional software risks. Appendix B specifically notes challenges pertinent to GenAI and cloud environments, including:

- Dependency on vast datasets which may harbor biases or quality issues, and are susceptible to poisoning attacks [17].
- Risks associated with using pre-trained models, which can "increase levels of statistical uncertainty and cause issues with bias management, scientific validity, and reproducibility" [17, p.38]. This is crucial in cloud settings where models might be sourced from third parties.
- Increased opacity and difficulty in predicting failure modes or emergent behaviors, complicating security validation [17]. This aligns with the widely recognized 'black box' problems of LLMs, encompassing their complexity, lack of transparency, and limited explainability [16]. Specific security concerns not fully addressed by traditional frameworks, such as "evasion, model extraction, membership inference, availability, or other machine learning attacks" [17, p.39], including adversarial vulnerabilities common in LLMs [16].
- Specific security concerns not fully addressed by traditional frameworks, such as "evasion, model extraction, membership inference, availability, or other machine learning attacks" [17, p.39].
- Risks associated with "third-party AI technologies, transfer learning, and off-label use," which are highly relevant when using GenAI models hosted or integrated via cloud services [17, p.39].

2. **Trustworthiness Characteristics:** The RMF emphasizes achieving trustworthy AI by balancing several characteristics [17]. For security, the most critical are:

- **Secure and Resilient:** AI systems should maintain "confidentiality, integrity, and availability" and be able to "withstand unexpected adverse events or unexpected changes" [17, p.15]. This includes protecting against data poisoning, adversarial examples, and model exfiltration – key threats for GenAI. The RMF notes applicability of existing standards like the NIST Cybersecurity Framework here [17, p.15].
- **Accountable and Transparent:** While distinct from security, transparency and accountability are vital for security incident analysis, understanding



vulnerabilities, and assigning responsibility, especially in complex cloud supply chains [17].

- **Privacy-Enhanced:** GenAI often processes vast amounts of data, potentially including sensitive information. Privacy risks are intertwined with security, as data breaches impact both. The RMF advocates for privacy considerations throughout the lifecycle and mentions Privacy-Enhancing Technologies[17].
- **Valid and Reliable:** Systems must perform accurately and consistently. Unreliable GenAI could produce insecure code, faulty security recommendations, or fail in ways that create security openings [17].

3. **Risk Management Core Functions:** The RMF outlines four functions to operationalize risk management:

- **Govern:** Establishing a risk management culture, policies, accountability structures, and processes. Crucially, this includes policies addressing risks from "third-party software and data and other supply chain issues", vital for cloud-based GenAI [17, pp.21-24].
- **Map:** Establishing context, categorizing the AI system, understanding capabilities and limitations, and mapping risks/benefits, explicitly including those from third-party components[17].
- **Measure:** Applying methods and metrics to assess risks and evaluate trustworthy characteristics, including specific evaluations for security and resilience and privacy[17].
- **Manage:** Prioritizing and responding to risks, including managing risks from third-party entities and implementing incident response and recovery plans[17].

In essence, the NIST AI RMF 1.0 provides a comprehensive framework that, while voluntary and high-level, guides organizations in systematically considering the multifaceted risks, including significant security and privacy challenges, inherent in developing, deploying, and using complex AI systems like GenAI, particularly within the context of third-party dependencies common in cloud computing environments. It stresses the importance of integrating risk management throughout the AI lifecycle and addressing the unique characteristics and vulnerabilities of AI technologies.

Complementing overarching frameworks like the NIST AI RMF, specific architectural approaches are emerging to address the unique security challenges of GenAI in cloud environments. One prominent example is Zero Trust Architecture (ZTA) [16]. ZTA moves away from traditional perimeter-based security towards a model where trust is never assumed, and verification is continuously required[16]. This aligns well with the NIST RMF's emphasis on secure and resilient systems and proactive risk management, particularly given the 'black box' nature and dynamic deployment of many GenAI models [16], [17]. Key tenets include strict identity verification, micro-segmentation to limit lateral movement, least privilege access control, and continuous monitoring [16]. Implementing ZTA for LLMs involves specific considerations such as unified identity management across cloud platforms, AI-driven dynamic access policies, automated network segmentation, robust data encryption and classification, continuous threat monitoring tailored to LLM vulnerabilities, and ensuring compliance [16]. Interestingly, AI itself can enhance ZTA through behavioral analytics for continuous authentication or threat intelligence processing[16]. However, implementing ZTA

effectively presents its own challenges, including complexity, integration with legacy systems, resource requirements, and potential performance impacts [16].

### Balance of Automation and Human Oversight

The integration of Artificial Intelligence (AI), particularly Generative AI (GenAI), into cybersecurity presents a significant paradigm shift, offering powerful automation capabilities to counter increasingly sophisticated cyber threats. A recurring theme in the literature, however, is the inherent tension between the compelling benefits derived from this automation and the indispensable necessity of human oversight [2]. While AI-powered security automation provides crucial safeguards against evolving cyber dangers, the unique characteristics and potential risks associated with AI systems, especially GenAI, underscore the continued importance of human expertise and intervention [2], [3].

A fundamental principle, strongly articulated within risk management frameworks, is that no "high-risk" AI system should be operated without substantial human oversight [17, p.7]. This necessitates careful deliberation regarding whether the potential benefits of deploying such systems truly outweigh the potential negative impacts and risks [17]. In cybersecurity contexts, high-risk applications might include automated incident response systems with the potential for disruptive countermeasures, security policy generation influencing critical infrastructure, or threat analysis tools whose outputs directly inform high-stakes decisions. The NIST AI Risk Management Framework (AI RMF) emphasizes that in situations where AI systems present unacceptable negative risk levels, such as imminent significant negative impacts or the occurrence of severe harms, their development and deployment should cease until these risks can be sufficiently managed [17].

Despite the promising applications of GenAI for security automation—such as generating security reports, suggesting code fixes, or creating configuration scripts significant challenges remain in striking the right balance between automation and appropriate human oversight. Research highlights several critical issues stemming from the use of GenAI in automated security operations [3]. One major concern is the potential for over-dependence on AI tools, which could lead to complacency or a degradation of human skills [3]. Furthermore, GenAI models themselves are susceptible to adversarial risks, including data poisoning or prompt injection attacks designed to manipulate their outputs, presenting unique security challenges [3]. The inherent complexity and often opaque nature of decision-making processes within sophisticated AI systems, including GenAI, can also hinder effective oversight and accountability [17] [3].

Effectively managing GenAI in cybersecurity demands a recognition that complete automation without human intervention introduces unacceptable risks [3]. Human oversight is crucial not merely as a final checkpoint but throughout the AI lifecycle. This includes defining system goals and constraints, interpreting ambiguous or novel situations that fall outside the AI's training data, providing contextual understanding that the AI may lack, and making ethical judgments, particularly when potential actions have significant consequences [17]. The NIST AI RMF emphasizes the importance of clearly defined human roles and responsibilities within human-AI configurations, acknowledging the influence of human cognitive biases and the need for systems that are explainable and interpretable to those operating or overseeing them [17].

Frameworks like the NIST AI RMF provide structured approaches to managing these challenges. The GOVERN function stresses establishing a risk management culture, defining roles, and ensuring accountability [17, p. 21-24]. The MAP function requires establishing context, understanding system limitations, and defining processes for human oversight [17, p. 24-28]. MEASURE involves ongoing monitoring of performance, safety, and fairness, incorporating feedback mechanisms [17, p. 28-31]. Crucially, the MANAGE function includes planning risk responses and implementing mechanisms to supersede, disengage, or deactivate AI systems demonstrating performance inconsistent with intended use, alongside robust post-deployment monitoring and incident response plans [17, p. 31-33].

Ultimately, the effective use of GenAI in cybersecurity hinges on achieving a balanced, symbiotic relationship between automated capabilities and human expertise. This balanced approach acknowledges the complementary strengths of humans and AI. GenAI can process vast amounts of data and automate repetitive tasks at scale and speed, while humans provide critical thinking, contextual awareness, ethical guidance, and ultimate accountability [3]. Preventive efforts and well-planned action plans, incorporating robust human oversight mechanisms, are essential to harness the benefits of GenAI for cybersecurity while mitigating its inherent risks [3].

### Summary Literature review

This literature review demonstrates that Generative AI (GenAI) represents a transformative technology for security automation within hyperscale cloud environments. The analysis reveals significant potential for GenAI to enhance security operations through automated threat detection, policy generation, and incident response, particularly across complex multi-cloud settings. Research highlights notable advancements in conceptual frameworks for multi-cloud policy orchestration, validation mechanisms to ensure trust and accuracy, and technical approaches for implementing GenAI at scale. The most promising strategies often leverage multi-cloud architectures, zero-trust principles, and comprehensive security frameworks, while necessarily acknowledging the unique infrastructure requirements of GenAI itself. However, despite this progress, persistent challenges related to trust, validation, data privacy and quality, and the crucial balance between automation and human oversight remain significant considerations. As this field continues its rapid evolution, interdisciplinary collaboration will be essential to develop robust ethical norms and innovative defense mechanisms, addressing current issues while guiding the responsible application of GenAI in cybersecurity.

## 2.5 Research Gaps



## **Appendix A**

# **Appendix Title Here**

Write your Appendix content here.



# Bibliography

- [1] K. Khanna, "ENHANCING CLOUD SECURITY WITH GENERATIVE AI: EMERGING STRATEGIES AND APPLICATIONS," *JARET*, vol. 3, no. 1, pp. 234–244, Jun. 14, 2024, Number: 1 Publisher: IAEME Publication, issn: 2295-5152. Accessed: Apr. 8, 2025. [Online]. Available: [https://iaeme.com/Home/article\\_id/JARET\\_03\\_01\\_021](https://iaeme.com/Home/article_id/JARET_03_01_021).
- [2] D. K. Seth, K. K. Ratra, and A. P. Sundareswaran, "AI and generative AI-driven automation for multi-cloud and hybrid cloud architectures: Enhancing security, performance, and operational efficiency," *2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC)*, pp. 00 784–00 793, Jan. 6, 2025, Conference Name: 2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC) ISBN: 9798331507695 Place: Las Vegas , NV, USA Publisher: IEEE. doi: [10.1109/CCWC62904.2025.10903928](https://doi.org/10.1109/CCWC62904.2025.10903928). Accessed: Apr. 8, 2025. [Online]. Available: <https://ieeexplore.ieee.org/document/10903928/>.
- [3] A. Patel, P. Pandey, H. Ragothaman, R. Molleti, and D. R. Peddinti, "Generative AI for automated security operations in cloud computing," *2025 IEEE 4th International Conference on AI in Cybersecurity (ICAIC)*, pp. 1–7, Feb. 5, 2025, Conference Name: 2025 IEEE 4th International Conference on AI in Cybersecurity (ICAIC) ISBN: 9798331518882 Place: Houston, TX, USA Publisher: IEEE. doi: [10.1109/ICAIC63015.2025.10849302](https://doi.org/10.1109/ICAIC63015.2025.10849302). Accessed: Mar. 31, 2025. [Online]. Available: <https://ieeexplore.ieee.org/document/10849302/>.
- [4] "Securing generative AI: Introduction to the generative AI security scoping matrix," Amazon Web Services, Inc. Accessed: Apr. 9, 2025. [Online]. Available: <https://aws.amazon.com/ai/generative-ai/security/scoping-matrix/>.
- [5] "Securing generative AI: An introduction to the generative AI security scoping matrix | AWS security blog." Section: Amazon Bedrock, Accessed: Apr. 8, 2025. [Online]. Available: <https://aws.amazon.com/blogs/security/securing-generative-ai-an-introduction-to-the-generative-ai-security-scoping-matrix/>.
- [6] C. Y. Haryanto, M. H. Vu, T. D. Nguyen, E. Lomempow, Y. Nurliana, and S. Taheri, *SecGenAI: Enhancing security of cloud-based generative AI applications within australian critical technologies of national interest*, Jul. 1, 2024. doi: [10.48550/arXiv.2407.01110](https://doi.org/10.48550/arXiv.2407.01110). arXiv: 2407.01110[cs]. Accessed: Aug. 26, 2024. [Online]. Available: <http://arxiv.org/abs/2407.01110>.
- [7] D. o. I. S. a. Resources. "Australia's AI ethics principles | australia's artificial intelligence ethics principles | department of industry science and resources," <https://www.industry.gov.au/node/91877>, Accessed: Apr. 9, 2025. [Online]. Available: <https://www.industry.gov.au/publications/australias-artificial-intelligence-ethics-principles/australias-ai-ethics-principles>.
- [8] "ISO/IEC 38500:2024," ISO, Accessed: Apr. 9, 2025. [Online]. Available: <https://www.iso.org/standard/81684.html>.

- [9] Sushil Prabhu Prabhakaran, "Integration patterns in unified AI and cloud platforms: A systematic review of process automation technologies," *Int. J. Sci. Res. Comput. Sci. Eng. Inf. Technol*, vol. 10, no. 6, pp. 1932–1940, Dec. 15, 2024, ISSN: 2456-3307. DOI: [10.32628/CSEIT241061229](https://doi.org/10.32628/CSEIT241061229). Accessed: Apr. 8, 2025. [Online]. Available: <https://ijsrcseit.com/index.php/home/article/view/CSEIT241061229>.
- [10] D. Bringhenti, R. Sisto, and F. Valenza, "Security automation for multi-cluster orchestration in kubernetes," *2023 IEEE 9th International Conference on Network Softwarization (NetSoft)*, pp. 480–485, Jun. 19, 2023, Conference Name: 2023 IEEE 9th International Conference on Network Softwarization (NetSoft) ISBN: 9798350399806 Place: Madrid, Spain Publisher: IEEE. DOI: [10.1109/NetSoft57336.2023.10175419](https://doi.org/10.1109/NetSoft57336.2023.10175419). Accessed: Apr. 8, 2025. [Online]. Available: <https://ieeexplore.ieee.org/document/10175419/>.
- [11] K. Hammar and R. Stadler, "Digital twins for security automation," *NOMS 2023-2023 IEEE/IFIP Network Operations and Management Symposium*, pp. 1–6, May 8, 2023, Conference Name: NOMS 2023-2023 IEEE/IFIP Network Operations and Management Symposium ISBN: 9781665477161 Place: Miami, FL, USA Publisher: IEEE. DOI: [10.1109/NOMS56928.2023.10154288](https://doi.org/10.1109/NOMS56928.2023.10154288). Accessed: Apr. 8, 2025. [Online]. Available: <https://ieeexplore.ieee.org/document/10154288/>.
- [12] Z. Zhou, C. Zhongwen, Z. Tiecheng, and G. Xiaohui, "The study on network intrusion detection system of snort," May 1, 2010. DOI: [10.1109/ICNDS.2010.5479341](https://doi.org/10.1109/ICNDS.2010.5479341).
- [13] S. Surathunmanun, W. Ongsakul, and J. G. Singh, "Exploring the role of generative artificial intelligence in the energy sector: A comprehensive literature review," *2024 International Conference on Sustainable Energy: Energy Transition and Net-Zero Climate Future (ICUE)*, pp. 1–11, Oct. 21, 2024, Conference Name: 2024 International Conference on Sustainable Energy: Energy Transition and Net-Zero Climate Future (ICUE) ISBN: 9798331517076 Place: Pattaya City, Thailand Publisher: IEEE. DOI: [10.1109/ICUE63019.2024.10795598](https://doi.org/10.1109/ICUE63019.2024.10795598). Accessed: Apr. 8, 2025. [Online]. Available: <https://ieeexplore.ieee.org/document/10795598/>.
- [14] "DevSecOps sentinel: GenAI-driven agentic workflows for comprehensive supply chain security | semantic scholar," Accessed: Apr. 8, 2025. [Online]. Available: <https://www.semanticscholar.org/paper/DevSecOps-Sentinel%3A-GenAI-Driven-Agentic-Workflows-Pillala-Azarpazhooh/c6936c7dcb49d540014eeb733bbacf>
- [15] R. L. V. Nyoto, M. Devega, and N. Nyoto, "Cyber security risks in the rapid development of generative artificial intelligence: A systematic literature review," *ComniTech : Journal of Computational Intelligence and Informatics*, vol. 1, no. 2, pp. 57–66, Dec. 29, 2024, ISSN: 3063-0630. Accessed: Apr. 8, 2025. [Online]. Available: <https://journal.unilak.ac.id/index.php/ComniTech/article/view/24539>.
- [16] "Zero-trust architecture (ZTA): Designing an AI-powered cloud security framework for LLMs' black box problems | semantic scholar," Accessed: Apr. 8, 2025. [Online]. Available: [https://www.semanticscholar.org/paper/Zero-Trust-Architecture-\(ZTA\)%3A-Designing-an-Cloud-Dash/ce5062561a36f15ec1cd203705736d7ab33](https://www.semanticscholar.org/paper/Zero-Trust-Architecture-(ZTA)%3A-Designing-an-Cloud-Dash/ce5062561a36f15ec1cd203705736d7ab33)



- [17] E. Tabassi, "Artificial intelligence risk management framework (AI RMF 1.0)," National Institute of Standards and Technology (U.S.), Gaithersburg, MD, NIST AI 100-1, Jan. 26, 2023, NIST AI 100–1. doi: 10.6028/NIST.AI.100-1. Accessed: Apr. 8, 2025. [Online]. Available: <http://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-1.pdf>.