



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

Declaration of Authorship	iii
Abstract	vii
Acknowledgements	ix
1 Introduction	1
1.1 Instead of an introduction	1
2 Background and Related Work	3
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 State of the Art	3
2.4.1 Methodology	4
2.4.2 AI-Driven Security Approaches	4
2.4.3 GenAI Security Frameworks	5
2.4.4 Approaches for Automated Cloud Security	8
2.4.5 Agent-Based Approaches	10
2.4.6 Security Risks	11
2.4.7 Balance of Automation and Human Oversight	16
2.4.8 Summary Literature State of the Art	17
2.5 Research Gaps	17
A Appendix Title Here	19
Bibliography	21

List of Figures

List of Tables

Listings

List of Abbreviations

ABAC	A tttribute- B ased A ccess C ontrol
ACSC	A ustralian C yber S ecurity C entre
AI	A rtificial I ntelligence
AI RMF	A rtificial I ntelligence R isk M anagement F ramework
APP	A ustralian P rivacy P inciples
CIA	C onfidentiality, I ntegrity, and A vailability
CSP	C loud S ervice P rovider
DTA	D igital T ransformation A gency
GenAI	G enerative A rtificial I ntelligence
LLM	L arge L anguage M odel
MLOps	M achine L earning O perations
MTTD	M ean T ime to D etect
MTTR	M ean T ime to R esolve
RAG	R etrieval- A ugmented G eneration
RPO	R ecovery P oint O bjectives
RTO	R ecovery T ime O bjectives
SOC	S ecurity O perations C enter
SRM	S hared R esponsibility M odel
WAF	W eb A pplication F irewalls
ZTA	Z ero T rust A 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 ⁻¹)
ω	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 State of the Art

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.

2.4.1 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.

2.4.2 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.

2.4.3 GenAI Security Frameworks

The rapid integration of Generative AI (GenAI) into various domains, including security automation within hyperscale cloud platforms, necessitates robust frameworks to govern its development, deployment, and operation securely. Unlike traditional software, GenAI systems introduce unique risks stemming from their data dependency, complexity, potential for emergent behaviors, and socio-technical nature [4]. These risks include prompt injection, data leakage through model inversion or training data extraction, adversarial attacks, generation of harmful or biased content, and insecure code generation [5, 6]. Consequently, a multi-faceted approach to security is required, encompassing foundational risk management principles, organizational governance structures, technical control implementation, and context-specific guidance. This subchapter reviews several key frameworks and guides that collectively address the challenge of securing GenAI systems.

A foundational element in managing AI risks is provided by the Artificial Intelligence Risk Management Framework (AI RMF 1.0) developed by the U.S. National Institute of Standards and Technology (NIST) [4]. This voluntary, non-sector-specific framework aims to help organizations manage AI risks and promote trustworthy and responsible AI development and use. It is structured around four core functions: GOVERN, MAP, MEASURE, and MANAGE. The GOVERN function is cross-cutting, establishing a risk management culture and processes throughout the organization. MAP involves contextualizing risks and understanding potential impacts. MEASURE employs qualitative and quantitative tools to analyze, assess, and track AI risks. MANAGE focuses on prioritizing and acting on risks based on the previous functions [4]. The NIST AI RMF emphasizes characteristics of trustworthy AI systems, including validity and reliability, safety, security and resilience, accountability and transparency, explainability and interpretability, privacy-enhancement, and fairness with harmful bias managed [4]. It acknowledges the challenges specific to AI risk management, such as difficulties in risk measurement (especially with third-party components and emergent risks), defining risk tolerance, prioritizing risks, and integrating AI risk management into broader enterprise strategies [4].

Building upon similar principles but offering a perspective from a major cloud provider, Google's Secure AI Framework (SAIF) presents a conceptual framework inspired by security best practices applied to software development, adapted for AI-specific risks [6]. SAIF proposes six core elements for secure AI systems[6]:

1. Expand strong security foundations to the AI ecosystem by applying and adapting existing controls.
2. Extend detection and response to include AI-specific threats and outputs.
3. Automate defenses using AI itself where appropriate, while keeping humans in the loop.

4. Harmonize platform-level controls to ensure consistency and prevent fragmentation.
5. Adapt controls with faster feedback loops, incorporating insights from red teaming and novel attack awareness.
6. Contextualize AI system risks within surrounding business processes, including model risk management and shared responsibility.

SAIF advocates a practical implementation approach involving understanding the specific use case, assembling a multidisciplinary team, providing an AI primer for all stakeholders, and applying the six core elements iteratively [6]. Like the NIST AI RMF, SAIF highlights the socio-technical nature of AI risks and the importance of context.

While NIST and Google provide overarching frameworks, securing GenAI, particularly Large Language Models (LLMs), requires specific organizational structures and practices. The LLM and Generative AI Security Center of Excellence (CoE) Guide from OWASP addresses this by outlining how to establish a dedicated CoE [7]. The primary objective of such a CoE is "to develop and enforce security policies and protocols for generative AI applications, facilitate cross-departmental collaboration to harness expertise from various fields, educate and train teams on the ethical and secure use of generative AI technologies, and serve as an advisory body for AI-related projects and initiatives within the organization" [7, p.4]. This guide emphasizes the necessity of a multidisciplinary team, bringing together expertise from Cybersecurity, AI/ML Development, IT Operations, Legal, Compliance, Ethics, Governance, Risk Management, Data Science, and various other user groups [7]. Establishing clear objectives, Key Performance Indicators, roles, and responsibilities is crucial for the CoE's success. The guide suggests a phased implementation (Planning, Integration, Operationalization, Evaluation) and highlights the importance of leveraging both internal and external expertise to address challenges like communication barriers, resistance to change, and skill gaps [7]. The CoE structure directly supports the GOVERN function described in the NIST AI RMF and aligns with SAIF's recommendation to assemble a cross-functional team.

To translate high-level governance and risk management principles into concrete verification steps, the OWASP LLM Applications Cybersecurity and Governance Checklist provides a practical, actionable tool [8]. Derived from the well-known OWASP Top 10 for Large Language Model Applications project, this checklist offers a structured approach to assessing the security posture of LLM-based systems. It covers a wide array of control areas critical for LLM security, extending beyond purely technical vulnerabilities to encompass essential governance aspects [8]. Key domains addressed include Input Validation and handling, Output Encoding and Handling, Access Control and Authorization, Data Privacy and Confidentiality, Model Training and Fine-tuning Security, API and Integration Security, robust Logging and Monitoring, Incident Response Planning, and overall Governance, Risk, and Compliance [8]. The checklist serves as a valuable resource for development teams, security assessors, and governance bodies to systematically identify potential weaknesses, guide the implementation of specific security controls, and perform gap analyses against established best practices [8]. In the context of broader frameworks, this checklist acts as a practical instrument for the MEASURE and MANAGE functions outlined in the NIST AI RMF [4], providing specific checks aligned with the risks highlighted by SAIF [6] and the technical controls advocated by SecGenAI [5]. It furnishes the Center of Excellence [8] with a concrete tool for enforcing security policies and protocols.

Regarding a specific GenAI architecture, the SecGenAI framework focuses on enhancing the security of cloud-based GenAI applications, particularly Retrieval-Augmented Generation (RAG) systems, within the context of Australian critical technologies [5]. SecGenAI adopts an end-to-end perspective covering Functional, Infrastructure, and Governance requirements. Functionally, it analyzes RAG components and identifies root causes for security concerns like injection attacks, data leakage, and model inversion [5]. It proposes specific countermeasures such as input validation, robust access controls, data protection techniques, and model security measures[5]. On the infrastructure side, SecGenAI details requirements for sandboxing, secure database connections, network segmentation, external attack prevention, and disaster recovery within a cloud environment[5]. Governance requirements emphasize alignment with Australian AI Ethics Principles and Privacy Principles (APP), advocating for fairness, accountability, traceability, data protection, regular audits, reliability, transparency, and compliance, structured using the ISO 38500 Evaluate-Direct-Monitor cycle [9]. SecGenAI explicitly incorporates the Shared Responsibility Model, detailing Cloud Service Provider and customer obligations for GenAI security[5]. This framework provides a detailed blueprint for securing a specific, common GenAI patterns by integrating technical and governance controls, reflecting a practical application of the principles found in NIST AI RMF and SAIF.

The successful implementation of these security frameworks inherently depends on the underlying platform architecture. The paper "Integration patterns in unified AI and cloud platforms" provides context by reviewing how AI, MLOps, workflow orchestration, and data processing converge in cloud-native environments[10]. It highlights the importance of MLOps frameworks for lifecycle management, workflow orchestration engines for process automation, and robust data processing systems as core components [10]. Security considerations must be embedded within these components and their integration patterns. For instance, securing data pipelines within MLOps, ensuring secure communication in federated learning setups, or implementing access controls within workflow orchestration are crucial [5, 6, 10]. The paper implicitly underscores that security cannot be an afterthought but must be integrated into the design and automation processes of these unified platforms, aligning with the principles of secure-by-design advocated by frameworks like SAIF and SecGenAI.

Facilitating this integration within a specific hyperscale cloud platform, the AWS GenAI Security Scoping Matrix serves as a practical aid for organizations utilizing AWS[11]. This matrix is designed explicitly to help customers navigate the complexities of the Shared Responsibility Model (SRM) as it applies to GenAI workloads deployed on AWS. It systematically maps common GenAI architectural components and layers against key security domains[11]. For each intersection of a GenAI component and a security domain, the matrix clarifies whether the responsibility for implementing controls lies primarily with AWS, the customer, or is shared between them [11]. This structured delineation is crucial for organizations to understand their specific security obligations when building and operating GenAI systems on AWS. Furthermore, the matrix guides customers in identifying and scoping the relevant AWS security services needed to fulfill their responsibilities [11]. It acts as a translator, converting the high-level principles of frameworks like NIST AI RMF [4] and SAIF [6], and the specific control requirements suggested by SecGenAI [5] or the OWASP Checklist [8], into actionable configurations and service selections within the AWS ecosystem. It directly supports the practical implementation of the SRM, a concept emphasized across multiple frameworks [4–6, 8], thereby enabling organizations to effectively manage risks within their AWS environment.

In summary, these frameworks and guides offer a layered approach to GenAI security. The NIST AI RMF provides a foundational, risk-based structure and defines trustworthiness. Google's SAIF offers a conceptual implementation path with core security elements, emphasizing adaptation and integration. The OWASP LLM CoE Guide focuses on the essential organizational and collaborative structures needed for effective governance. SecGenAI provides a detailed, integrated blueprint for securing a specific architecture, demonstrating how broader principles can be applied in context, including alignment with regional regulations. Practical tools like the OWASP LLM Checklist and provider-specific resources like the AWS Scoping Matrix aid in the MEASURE and MANAGE functions by providing concrete checks and configurations. The insights from [10] remind us that these security measures must be seamlessly integrated within the complex fabric of unified AI and cloud platforms, particularly within MLOps and automation workflows, to be effective. The recurring theme of the Shared Responsibility Model across multiple frameworks [4–6, 8] highlights the collaborative nature of securing GenAI in cloud environments. Collectively, these resources provide a comprehensive toolkit for organizations aiming to leverage GenAI for security automation and other critical tasks on hyperscale cloud platforms, enabling them to manage risks effectively and build trustworthy AI systems.

2.4.4 Approaches for Automated Cloud Security

This subsection details specific technical approaches and architectural patterns crucial for enabling automated cloud security. It reviews research on unified platforms integrating AI and MLOps across multi-cloud environments, techniques for automated policy orchestration in complex Kubernetes setups, and the application of digital twins for robust security policy validation. These approaches represent concrete mechanisms for realizing the potential of automation in dynamic cloud infrastructures.

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[12]. 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[12]. 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[12].

The proposed solution involves a top-level entity named the "Multi-Cluster Orchestrator"[12]. This orchestrator acts as a central management point, receiving inputs from managers of different domains[12]. These inputs include:

- A description of each domain's structure[12].
- High-level security requirements specifying allowed communications[12]. These requirements can be defined using an extended YAML syntax with special labels that abstract away low-level details[12].

Based on these inputs, the Multi-Cluster Orchestrator refines the high-level requirements into concrete configurations through a two-step process[12]:

1. It generates a "Global Configuration" that tracks communication pairs between services and required links between clusters, optimizing the overall cluster mesh setup[12].
2. It derives "Single Configurations" for each individual cluster, containing the specific parameters needed to connect the cluster to the 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[12].

The implementation, known as Multi-Cluster Orchestrator, 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[12]. 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[12].

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[13]. This approach utilizes an emulation system specifically designed to create high-fidelity digital replicas of target IT infrastructures[13]. 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[13]. This capability avoids impacting operational workflows on the real-world infrastructure[13].

The digital twin approach, following the research by Hammar and Stadler, enables a closed-loop learning process for crafting and refining security policies[13]. It starts with generating a digital twin of the target infrastructure. This is achieved using an emulation system constructed with virtualization tools like Docker containers, alongside virtual links and switches[13]. Within this digital twin, various security scenarios involving emulated attackers, defenders, and client populations are executed[13]. During these runs, monitoring agents collect detailed system measurements and logs, channeling this data through pipelines for analysis[13]. The gathered data and statistics are then used to build simulations[13]. Reinforcement learning techniques are applied to these simulations to learn potentially optimal security policies[13]. Finally, the performance of these learned policies is rigorously evaluated back in the digital twin, allowing for validation and further iteration[13].

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[13]. 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[14]. 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[13].

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.[5] 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.[5]

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 [15]. 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 [15].

2.4.5 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[16], 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[16].

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[16]. 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[16].

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

- **Automated Vulnerability and Impact Analysis:** Leveraging GenAI to analyze code, dependencies and infrastructure configurations for potential threats, assessing their potential impact in context[16].
- **Adaptive Compliance and Release Gating:** Enforcing security policies and compliance requirements dynamically, acting as automated checks before deployment[16].

- **Predictive Security:** Utilizing AI to identify potential future risks based on historical data and emerging threat patterns[16].

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[16].
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[16]. 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[16].

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.

2.4.6 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. [17], analyzing 17 relevant studies according to PRISMA 2020 guidelines [18], identifies several significant cybersecurity threats stemming primarily from the irresponsible application of GenAI technology. Complementing this, Surathunmanun et al. [15], 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 [17]. Furthermore, GenAI facilitates the creation of effective ransomware and diverse forms of malware, potentially empowering individuals with limited coding expertise to launch attacks [17]. Beyond typical malware, it can also generate other executable attack code, such as SQL injection scripts [17]. This potential for misuse is a major concern, where uncontrolled access or

improper application can lead to significant harm [15]. This includes leveraging GenAI to bypass security controls through techniques like prompt injection or jailbreaking, as highlighted in literature concerning Large Language Models. [15].

- **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 [17]. Concurrently, GenAI models are prone to generating plausible but factually incorrect or nonsensical information, known as hallucinations [15, 17]. This issue often arises from poor quality training data or suboptimal parameter settings [15] and can be exacerbated by data poisoning during the model training phase [17]. In a security context, hallucinations can manifest as faulty threat analyses, incorrect vulnerability assessments, or misleading security recommendations [15]. 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 [15]. 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 [15]. 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 [19].
- **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 [17]. User interactions and prompts can also be incorporated into training data, leading to potential data leakage and privacy violations [17]. This raises substantial intellectual property concerns and challenges compliance with regulations like GDPR [17]. The lack of transparency and control over how data is utilized presents considerable privacy risks [17]. Furthermore, insecure data handling practices can create security vulnerabilities [15]. Specific risks associated with LLMs, often used in cloud-hosted GenAI services, include inference attacks, data extraction attacks, data poisoning, supply chain vulnerabilities [15] and vulnerabilities to adversarial attacks stemming from the models complex and often opaque nature [19].
- **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 [15]. In automated security workflows operating in cloud environments, such errors could propagate rapidly, leading to service disruptions, critical misconfigurations, or a failure to respond effectively to genuine threats [15].

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 [15] and architectures like Zero Trust [15, 19], establishing clear ethical guidelines, and potentially developing new regulations specific to GenAI development and deployment [17]. 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 [15].

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) [4] 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" [4, 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 [4].

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 [4].
- Risks associated with using pre-trained models, which can "increase levels of statistical uncertainty and cause issues with bias management, scientific validity, and reproducibility" [4, 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 [4]. This aligns with the widely recognized 'black box' problems of LLMs, encompassing their complexity, lack of transparency, and limited explainability [19]. Specific security concerns not fully addressed by traditional frameworks, such as "evasion, model extraction, membership inference, availability, or other machine learning attacks" [4, p.39], including adversarial vulnerabilities common in LLMs [19].
- Specific security concerns not fully addressed by traditional frameworks, such as "evasion, model extraction, membership inference, availability, or other machine learning attacks" [4, 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 [4, p.39].

2. **Trustworthiness Characteristics:** The RMF emphasizes achieving trustworthy AI by balancing several characteristics [4]. 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" [4, 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 [4, 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 [4].

- **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[4].
- **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 [4].

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 [4, 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[4].
- **Measure:** Applying methods and metrics to assess risks and evaluate trustworthy characteristics, including specific evaluations for security and resilience and privacy[4].
- **Manage:** Prioritizing and responding to risks, including managing risks from third-party entities and implementing incident response and recovery plans[4].

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.

Adding to 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) [19]. ZTA moves away from traditional perimeter-based security towards a model where trust is never assumed, and verification is continuously required[19]. 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 [4, 19]. Key tenets include strict identity verification, micro-segmentation to limit lateral movement, least privilege access control, and continuous monitoring [19]. 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 [19]. Interestingly, AI itself can enhance ZTA through behavioral analytics for continuous authentication or threat intelligence processing[19]. However, implementing ZTA effectively presents its own

challenges, including complexity, integration with legacy systems, resource requirements, and potential performance impacts [19].

2.4.7 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 [4, p.7]. This necessitates careful deliberation regarding whether the potential benefits of deploying such systems truly outweigh the potential negative impacts and risks [4]. 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 [4].

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 [4] [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 [4]. 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 [4].

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 [4, p. 21-24]. The MAP function requires establishing context, understanding system limitations, and defining processes for human oversight [4, p. 24-28]. MEASURE involves ongoing monitoring of performance, safety, and fairness, incorporating feedback mechanisms [4, 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 [4, 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].

2.4.8 Summary Literature State of the Art

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.
- [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] 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](https://doi.org/10.6028/NIST.AI.100-1). Accessed: Apr. 8, 2025. [Online]. Available: <http://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-1.pdf>.
- [5] 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\[cs\]](https://arxiv.org/abs/2407.01110). Accessed: Aug. 26, 2024. [Online]. Available: <http://arxiv.org/abs/2407.01110>.
- [6] R. Hansen and P. Venables. "Introducing google's secure AI framework," Google, Accessed: Apr. 25, 2025. [Online]. Available: <https://blog.google/technology/safety-security/introducing-googles-secure-ai-framework/>.
- [7] OWASP Top 10 for LLMs Project, "LLM and GenAI Security Center of Excellence Guide," OWASP Foundation, Guide, version 1.0, Oct. 1, 2024, First Release Candidate based on revision history. Part of the OWASP Top 10 for Large Language Model Applications project. Accessed: Apr. 28, 2025. [Online]. Available: <https://ppl-ai-file-upload.s3.amazonaws.com/web/direct-files/attachments/22117806/aa82a187-e3a3-4834-a520-136f4f8667c6/OWASP-Top-10-for-LLM-COE-Guide-v1.0-1.pdf>.

- [8] OWASP Top 10 for LLM Applications Project, "LLM AI Cybersecurity & Governance Checklist," OWASP Foundation, Checklist, version 1.1, Apr. 10, 2024, Based on revision history. Contributors listed include Sandy Dunn. Accessed: Apr. 28, 2025. [Online]. Available: https://ppl-ai-file-upload.s3.amazonaws.com/web/direct-files/attachments/22117806/c36e9109-6614-49ff-b707-2d580de2a140/LLM_AI_Security_and_Governance_Checklist-v1.1.pdf.
- [9] "ISO/IEC 38500:2024," ISO, Accessed: Apr. 9, 2025. [Online]. Available: <https://www.iso.org/standard/81684.html>.
- [10] 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>.
- [11] "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/>.
- [12] 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/>.
- [13] 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/>.
- [14] 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).
- [15] 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/>.

- [16] G. Pillala, D. Azarpazhooh, and S. Baxter, "DevSecOps sentinel: GenAI-driven agentic workflows for comprehensive supply chain security," *Computer and Information Science*, vol. 18, no. 1, p39, Dec. 20, 2024, Number: 1, ISSN: 1913-8989. doi: [10.5539/cis.v18n1p39](https://ccsenet.org/journal/index.php/cis/article/view/0/51118). Accessed: Apr. 27, 2025. [Online]. Available: <https://ccsenet.org/journal/index.php/cis/article/view/0/51118>.
- [17] 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>.
- [18] M. J. Page et al., "The PRISMA 2020 statement: An updated guideline for reporting systematic reviews," *BMJ*, n71, Mar. 29, 2021, ISSN: 1756-1833. doi: [10.1136/bmj.n71](https://www.bmj.com/lookup/doi/10.1136/bmj.n71). Accessed: Apr. 12, 2025. [Online]. Available: <https://www.bmj.com/lookup/doi/10.1136/bmj.n71>.
- [19] B. Dash, *Zero-trust architecture (ZTA): Designing an AI-powered cloud security framework for LLMs' black box problems*, Rochester, NY, Mar. 12, 2024. doi: [10.2139/ssrn.4726625](https://papers.ssrn.com/abstract=4726625). Accessed: Apr. 27, 2025. [Online]. Available: <https://papers.ssrn.com/abstract=4726625>.