



(REVIEW ARTICLE)



AI-driven anomaly detection in cloud computing environments

Chukwuemeka Nwachukwu ^{1,*}, Kehinde Durodola-Tunde ² and Chukwuebuka Akwiwu-Uzoma ³

¹ Department of AI and Research, University of Bradford, UK.

² Department of AI, Product Management, MoniMoore, London, UK.

³ Department of Computing, University of Dundee, Dundee, United Kingdom.

International Journal of Science and Research Archive, 2024, 13(02), 692–710

Publication history: Received on 03 October 2024; revised on 09 November 2024; accepted on 11 November 2024

Article DOI: <https://doi.org/10.30574/ijrsra.2024.13.2.2184>

Abstract

The rapid adoption of cloud computing has changed the way businesses manage and store data, but it has also introduced new security challenges. One of the most pressing concerns in cloud environments is the detection of anomalies, which can signal potential security breaches, system failures, or performance issues. Traditional anomaly detection methods often fall short due to the complexity, scalability, and dynamic nature of cloud infrastructures. In recent years, Artificial Intelligence (AI)-driven anomaly detection techniques, particularly those leveraging machine learning and deep learning, have shown promise in overcoming these limitations. This paper reviews AI-driven approaches to anomaly detection in cloud computing environments, exploring their applications in enhancing cloud security, optimizing performance, and ensuring efficient resource management. The paper examines the strengths of various AI techniques, including supervised and unsupervised learning, deep learning, and hybrid models, highlighting their capacity to detect complex, previously unknown anomalies. Despite their advantages, implementing AI-based systems in cloud environments presents challenges, including data quality issues, scalability concerns, and computational resource requirements. Solutions such as federated learning and model optimization techniques are explored as methods to address these challenges. Furthermore, the paper discusses future research directions, including the integration of AI-driven anomaly detection with emerging technologies like blockchain and IoT, and the potential for advancements in self-supervised learning and explainable AI (XAI). This review concludes by emphasizing the critical role of AI in securing cloud infrastructures and the promising future of anomaly detection in the cloud computing landscape.

Keywords: Cloud Computing; Anomaly Detection; Artificial Intelligence; Machine Learning; Deep Learning; Security

1. Introduction

1.1. Context and Importance

Cloud computing has become an integral part of modern IT infrastructures, enabling organizations to scale operations, enhance efficiency, and reduce costs. The adoption of cloud-based services, including storage, computing power, and networking, has surged in recent years, as businesses migrate away from traditional on-premise solutions. This shift is largely driven by the cloud's flexibility, scalability, and pay-as-you-go model, which allows companies to dynamically adjust resources based on demand. Major cloud providers, such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform, offer a range of solutions, fostering a robust ecosystem of cloud services used by a diverse set of industries.

However, the growing reliance on cloud environments has led to increased challenges, particularly in the areas of security and anomaly detection. As cloud environments become more complex, ensuring the integrity of data and

* Corresponding author: Chukwuemeka Nwachukwu

maintaining secure operations is a significant concern. Malicious attacks, data breaches, and system failures are persistent threats that require vigilant monitoring and rapid response. Traditional security methods, such as firewalls and intrusion detection systems, are often insufficient for the dynamic and large-scale nature of cloud environments. This has created a pressing need for advanced techniques capable of identifying abnormal behaviour and potential threats in real-time (Chandola et al., 2009).

AI, specifically machine learning (ML), holds great potential in addressing these challenges. AI-driven anomaly detection can continuously monitor cloud infrastructures, adapt to new threats, and provide a level of insight that traditional methods cannot match. This approach enables more efficient and effective responses to security risks, enhancing the overall security posture of cloud systems.

1.2. Anomaly Detection Overview

Anomaly detection is the process of identifying patterns in data that do not conform to expected behaviour. In cloud computing, anomaly detection refers to the ability to identify unusual activities, such as abnormal system performance or unauthorized access, that could indicate security breaches, fraud, or system malfunctions. These anomalies are often the precursors to more significant issues, including cyber-attacks or service disruptions, making early detection crucial in preventing damage (Iglewicz & Hoaglin, 1993).

Traditional anomaly detection methods, such as statistical models or rule-based systems, typically rely on predefined thresholds or historical data patterns to flag unusual behaviour. While these approaches can be effective for simple environments, they struggle to keep up with the complexity and dynamism of modern cloud infrastructures. Cloud environments feature ever-changing workloads, diverse data types, and numerous variables, making it difficult for traditional methods to adapt quickly to new patterns or evolving threats.

AI-driven anomaly detection, particularly through ML, offers a more sophisticated and dynamic approach. ML algorithms can analyse vast amounts of data from cloud systems, identify patterns, and learn from the data to detect deviations from normal behaviour without predefined rules (Chandola et al., 2009). These algorithms continuously improve by recognizing emerging patterns and adapting to new data, making them far more effective in detecting subtle anomalies that might go unnoticed by traditional methods. Additionally, AI-powered systems can process and analyse large-scale cloud data in real-time, enabling quicker identification and response to threats.

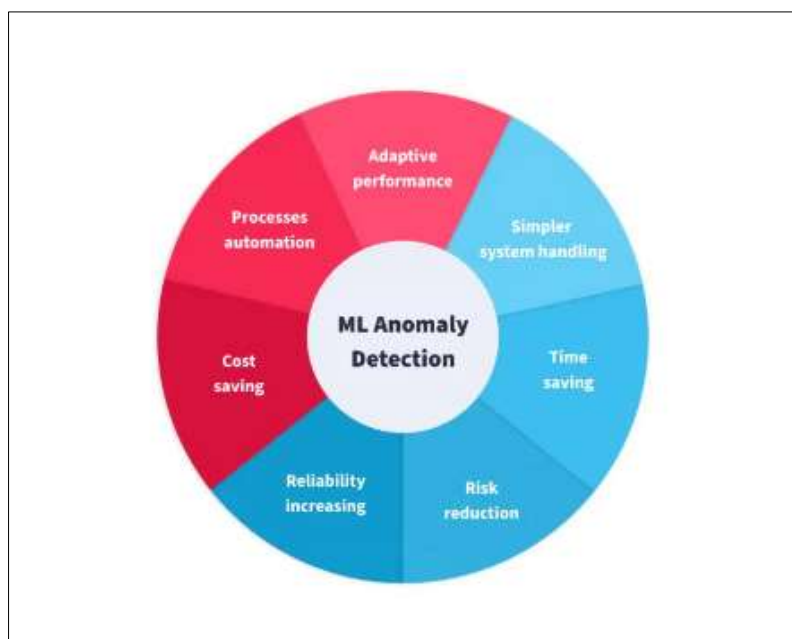


Figure 1 Significance of ML in Anomaly Detection (Chandola et al., 2009)

ML-based anomaly detection models, such as supervised and unsupervised learning algorithms, clustering techniques, and neural networks, have proven highly effective in cloud security contexts, reducing false positives and increasing detection accuracy (Chandola et al., 2009). These advanced methods allow cloud providers and enterprises to proactively address security concerns, ensuring the integrity of data and maintaining operational efficiency.

1.3. Aim of the Paper

The primary aim of this paper is to review and critically assess AI-driven anomaly detection approaches within cloud computing environments. As cloud systems continue to grow in scale and complexity, traditional anomaly detection methods often fall short in identifying subtle, sophisticated, or novel security threats. This paper seeks to bridge the knowledge gap by exploring how AI, particularly ML and deep learning techniques, enhances anomaly detection capabilities, allowing for more dynamic, accurate, and real-time identification of potential threats and irregularities. By evaluating various AI-driven strategies, including supervised and unsupervised learning models, clustering algorithms, and neural network architectures, this paper provides a comprehensive overview of how these advanced techniques outperform conventional methods in addressing cloud security concerns.

2. Cloud computing environments and security challenges

2.1. Cloud Computing Models and Architectures

Cloud computing encompasses various service models and deployment architectures that cater to different business and operational needs. The primary service models include Infrastructure as a Service (IaaS), Platform as a Service (PaaS), and Software as a Service (SaaS). IaaS offers virtualized computing resources over the internet, providing users with storage, servers, and networking components without the need for on-premise hardware. This model is ideal for businesses requiring flexibility and scalability. PaaS delivers an environment where developers can build, test, and deploy applications without managing the underlying infrastructure. This model simplifies the development process, allowing organizations to focus on coding and innovation. SaaS provides fully functional software applications hosted on the cloud, which users can access through web browsers. It offers convenience and cost efficiency, removing the need for installations or maintenance (Mell & Grance, 2011).

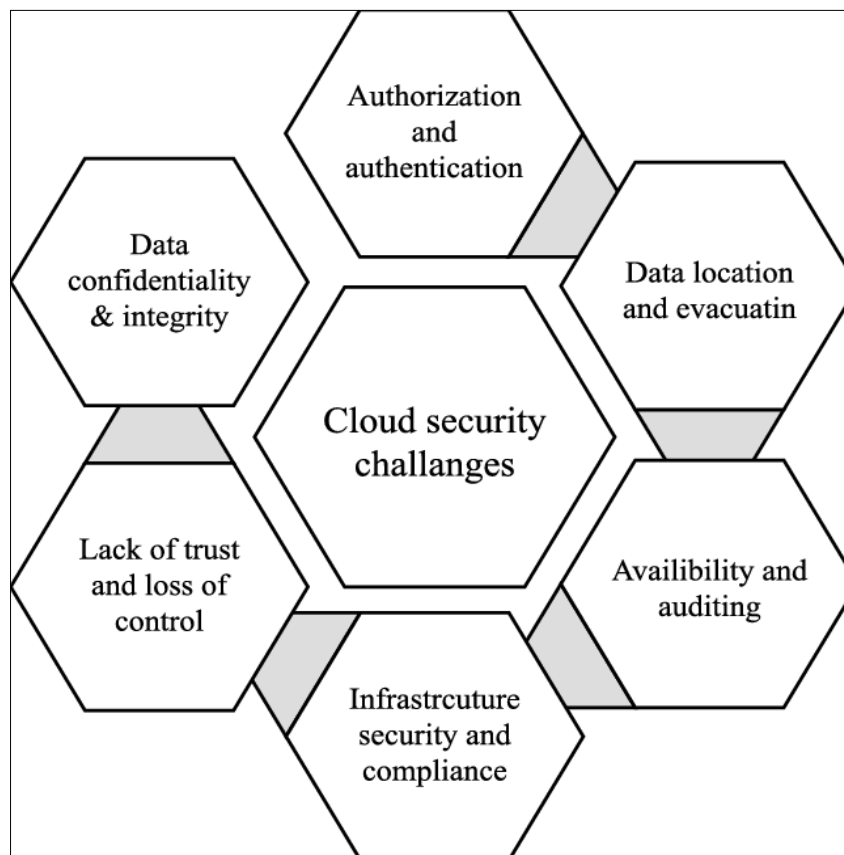


Figure 2 Cloud Security Challenge (Buyya et al., 2010).

Deployment architectures further define how cloud resources are utilized. Public clouds are shared infrastructures provided by third-party vendors and are cost-effective for many users. Private clouds are dedicated to a single organization, offering enhanced control and security. Hybrid clouds combine elements of both public and private

models, allowing data and applications to move between the two, offering a balance between scalability and privacy (Buyya et al., 2010).

The diverse and interconnected nature of these models introduces complexity, making it challenging to monitor and secure cloud environments. This complexity underscores the need for robust anomaly detection systems capable of handling diverse data sources and identifying potential threats across different cloud structures.

2.2. Security Risks and Threats in Cloud Environments

Cloud environments face a myriad of security risks and threats that can compromise data integrity, privacy, and service continuity. One significant concern is data breaches, where sensitive information is accessed by unauthorized entities. Such breaches can occur due to vulnerabilities in application programming interfaces (APIs), misconfigurations, or weak authentication measures. The resulting exposure of data can lead to financial losses, regulatory penalties, and damage to an organization's reputation (Zhang et al., 2010).

Distributed Denial of Service (DDoS) attacks are another common threat. These attacks aim to disrupt cloud services by overwhelming them with a flood of traffic, causing service degradation or outages. Cloud providers invest heavily in mitigating such attacks, but the scale and frequency of DDoS incidents make them an ongoing challenge (Mirkovic & Reiher, 2004). The dynamic and expansive nature of cloud systems can make it difficult to predict or prevent these attacks, necessitating continuous monitoring and anomaly detection to identify abnormal traffic patterns and respond swiftly.

Insider threats are particularly insidious, as they involve individuals with authorized access exploiting their position to steal or damage data. These threats can come from disgruntled employees, contractors, or third-party service providers. Insider activities often mimic legitimate user behaviour, making them hard to detect through traditional monitoring methods. Advanced anomaly detection systems that use ML can help identify subtle deviations in user behaviour that may indicate malicious intent (Greitzer & Frincke, 2010).

Unauthorized access due to weak authentication practices or credential theft is another major concern. Cloud services that are not configured with multi-factor authentication (MFA) or robust identity management protocols are vulnerable to unauthorized access. Attackers who gain entry through compromised credentials can escalate privileges and move laterally within the network, exfiltrating data or deploying malware. Implementing AI-driven anomaly detection can enhance security by monitoring access patterns and flagging atypical logins or behaviour that may signal a breach (Cai et al., 2016).

Other risks include misconfigurations of cloud resources, often due to human error or insufficient security policies. Such misconfigurations can expose data and systems to attackers. The prevalence of automated tools and infrastructure as code (IaC) increases the risk of deploying environments with incorrect security settings. Detecting these vulnerabilities promptly is crucial to minimizing potential damage.

The challenges in protecting cloud environments are amplified by their interconnected and often dynamic nature. The introduction of AI and ML in security protocols, particularly for anomaly detection, offers a significant advantage in addressing these issues. These technologies can process vast volumes of data, learn from past incidents, and detect potential threats in real-time by identifying deviations from established baselines.

2.3. Challenges in Detecting Anomalies in Cloud Computing

Implementing effective anomaly detection in cloud computing environments poses significant challenges due to the unique characteristics of these systems. Scalability is a primary concern, as cloud infrastructures often consist of numerous distributed resources that must be monitored simultaneously. Traditional detection systems struggle to process and analyse the vast amounts of data generated in real-time by these large-scale systems, leading to potential blind spots in security monitoring (Chandola et al., 2009).

The complexity of cloud environments also poses a hurdle. Cloud systems typically involve multiple layers, such as virtual machines, containers, and various services interconnected across regions and data centres. This complexity increases the potential for subtle, multifaceted anomalies that traditional rule-based or signature-based detection methods may miss (Zhang et al., 2010). The dynamic nature of cloud infrastructures further complicates detection efforts. Continuous changes, including scaling, shifting workloads, and frequent updates, can make it difficult to establish a consistent baseline for normal behaviour (Mell & Grance, 2011).

These challenges highlight the need for more adaptive and intelligent approaches to anomaly detection. Traditional methods, which rely on static rules or thresholds, lack the flexibility to adapt to evolving patterns in real-time. AI-driven methods, particularly those leveraging ML algorithms, can learn from data patterns, adjust to new behaviours, and scale efficiently with cloud environments (Cai et al., 2016). This allows them to address scalability and complexity issues, making them more effective in detecting anomalies that traditional approaches might overlook.

3. Anomaly detection techniques in cloud computing

3.1. Traditional Methods of Anomaly Detection

Classical anomaly detection techniques have long been used to identify irregularities in various systems. Among the most common approaches are statistical analysis, rule-based methods, and clustering algorithms. Each has been employed in different contexts to flag data points that deviate from the norm.

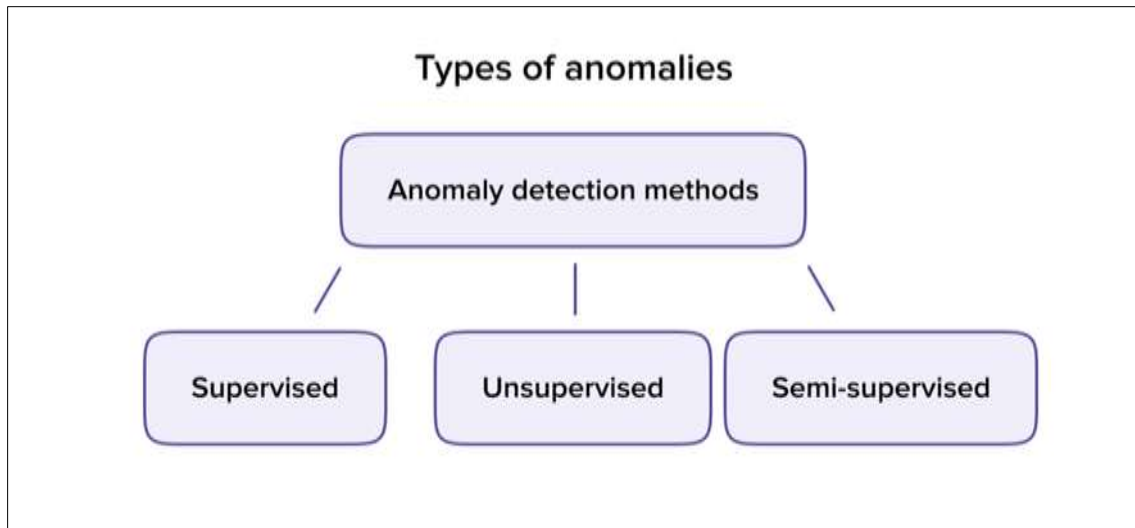


Figure 3 Types of Anomalies (Chandola et al., 2009).

- Statistical analysis techniques, such as z-score and probability distribution models, rely on historical data to identify outliers. These methods are effective for datasets with a consistent distribution where thresholds for "normal" behaviour can be established (Chandola et al., 2009). However, their effectiveness diminishes when applied to cloud environments, which often feature dynamic, non-stationary data streams. These environments require constant adjustment of baselines, making static statistical methods prone to false positives and negatives.
- Rule-based methods involve creating sets of predefined rules or conditions under which an anomaly is flagged. These methods are straightforward to implement and interpret. For example, if network traffic exceeds a certain limit or user activity deviates from an expected pattern, a rule-based system would trigger an alert (Iglewicz & Hoaglin, 1993). While effective for specific use cases, rule-based systems are inflexible and unable to adapt to new or unknown threats. As cloud systems evolve and threats become more sophisticated, maintaining and updating these rule sets becomes a cumbersome task that lacks scalability.
- Clustering algorithms, such as k-means and DBSCAN, group data points based on their similarity and can flag data points that do not belong to any cluster as potential anomalies. While clustering can identify irregular patterns without prior labels, it assumes that normal data points form well-defined clusters. In complex and highly variable cloud environments, data distributions may overlap, making it difficult for clustering algorithms to distinguish between legitimate variations and anomalies (Xu & Tian, 2015). Additionally, clustering methods often require parameter tuning, which may not be practical for rapidly changing cloud conditions.

These traditional methods share common limitations when applied to cloud environments. They struggle with scalability, as they are not designed to handle the volume and velocity of data generated by modern cloud infrastructures. Moreover, traditional methods often fail to adapt to the inherent dynamism of cloud systems, where baseline behaviours can shift due to scaling, load balancing, or continuous updates. AI-based anomaly detection methods overcome these limitations by leveraging ML algorithms that can handle large, high-dimensional datasets and

adapt to evolving patterns in real-time. These advanced methods provide more accurate and efficient solutions, enabling comprehensive monitoring of complex cloud environments.

3.2. AI-Based Anomaly Detection Techniques

AI-based techniques have revolutionized anomaly detection in cloud computing environments, offering improved scalability, adaptability, and precision. Key methods include supervised learning, unsupervised learning, deep learning, and reinforcement learning (RL). Each of these approaches contributes unique strengths to detecting irregularities in complex, high-dimensional data.

3.2.1. Supervised Learning

Supervised learning involves training models using labelled data, where each data point is annotated as either normal or anomalous. Common algorithms include support vector machines, decision trees, and random forests (Chukwunweike JN et al., 2024). These models learn to classify new data points based on patterns identified in the training set. Supervised learning techniques are highly effective when large, high-quality labelled datasets are available. In cloud computing, these models can be trained to detect known threats, such as unusual login attempts or abnormal network traffic (Hodge & Austin, 2004).

However, supervised learning has limitations in cloud anomaly detection. The dynamic nature of cloud environments often introduces new, unknown anomalies that the model has not been trained to recognize. Additionally, creating a comprehensive labelled dataset can be resource-intensive and time-consuming. This limitation underscores the need for models that can adapt to unseen data (Chandola et al., 2009).

3.2.2. Unsupervised Learning

Unsupervised learning is particularly advantageous for anomaly detection in cloud environments where labelled data is scarce or unavailable. These algorithms identify patterns and group data points based on their intrinsic characteristics (Adetunji A et al., 2024). Clustering algorithms, such as k-means and DBSCAN, are common unsupervised techniques that help identify outliers by detecting data points that do not conform to established groupings (Xu & Tian, 2015). In cloud security, clustering can reveal anomalies in network traffic or user behaviour that may indicate security breaches or misuse.

3.2.3 Principal Component Analysis (PCA) and Isolation Forests are also popular in unsupervised anomaly detection. PCA reduces the dimensionality of data, highlighting patterns that differ from the norm, while Isolation Forests isolate data points by randomly partitioning the dataset to highlight anomalies. These methods excel in cloud environments due to their ability to process large amounts of data without prior labelling (Liu et al., 2008).

The primary limitation of unsupervised learning is the potential for higher false positive rates, as the algorithm may flag legitimate variations as anomalies. Additionally, these methods can struggle with highly dynamic data streams that characterize cloud environments.

3.2.3. Deep Learning

Deep learning has become an integral part of modern anomaly detection due to its ability to capture complex, non-linear relationships in large datasets. Neural networks can process vast amounts of high-dimensional data, making them ideal for cloud computing applications (Chukwunweike et al., 2024). One of the most widely used deep learning techniques is the autoencoder. Autoencoders are neural network models that learn to encode data into a compressed format and then decode it back to its original form. Anomalies are detected when the reconstruction error (the difference between input and output) exceeds a certain threshold, indicating that the model struggles to accurately represent atypical data (Hinton & Salakhutdinov, 2006).

3.2.5 Recurrent Neural Networks (RNNs) and their advanced variants, such as Long Short-Term Memory (LSTM) networks, are also used in anomaly detection. These models excel at processing sequential data, making them suitable for detecting temporal anomalies in cloud environments, such as irregular user sessions or traffic patterns. By maintaining a memory of previous inputs, RNNs can identify context-dependent anomalies that traditional methods might overlook (Hochreiter & Schmidhuber, 1997). The power of deep learning lies in its adaptability and ability to generalize complex patterns without human intervention. However, deep learning models often require significant computational resources and can be difficult to interpret, posing challenges for their deployment in resource-constrained cloud settings.

3.3. RL

RL introduces a different approach to anomaly detection by training models to make decisions based on feedback from the environment. In the context of cloud anomaly detection, RL-based systems can learn adaptive strategies to respond to evolving threats (Chukwunweike JN et al...2024). These models can monitor real-time data and adjust their detection parameters dynamically to maintain high accuracy (Sutton & Barto, 2018). A notable application of RL in cloud security involves adaptive thresholding, where the model learns to modify the sensitivity of detection mechanisms based on the type and frequency of incoming data. RL can also optimize multi-step responses, enabling anomaly detection systems to not only identify anomalies but also recommend or execute corrective actions. This adaptability makes RL especially effective in environments where the threat landscape changes frequently.

Despite its advantages, RL requires substantial training and may not perform well in scenarios where immediate feedback is limited. Training RL models can also be complex, involving a careful balance between exploration (testing new strategies) and exploitation (refining known strategies). The increasing sophistication of AI methods in anomaly detection, driven by their capacity to process large datasets, adapt to changing conditions, and identify complex patterns, marks a significant leap forward in cloud security. By incorporating these advanced AI techniques, organizations can enhance their ability to detect anomalies and respond to potential threats in real-time, outperforming traditional approaches and bolstering the resilience of cloud infrastructures.

3.3.1. Hybrid Approaches

Hybrid models, which integrate multiple AI techniques, offer enhanced anomaly detection performance by combining the strengths of various methods. For instance, ensemble learning can merge supervised and unsupervised algorithms to detect both known and unknown anomalies (Chukwunweike JN et al...2024). Hybrid deep learning models, such as those combining autoencoders with clustering algorithms, enhance the detection of complex patterns while minimizing false positives. These approaches improve scalability and adaptability in cloud environments, addressing limitations inherent in single-method systems (Zhou et al., 2012). With an understanding of these advanced AI techniques, we can now explore their practical applications and how they are deployed to safeguard cloud infrastructures.

4. Practical applications of ai-driven anomaly detection in cloud environments

4.1. Application in Cloud Security

AI-driven anomaly detection is pivotal in strengthening cloud security by identifying and mitigating potential threats such as unauthorized access, DDoS attacks, and insider threats. The deployment of AI methods has enabled cloud security frameworks to become more proactive, adaptive, and accurate in identifying irregularities within vast data streams.

- Unauthorized access is a significant risk in cloud environments where sensitive data is stored and managed. AI-based anomaly detection systems can monitor login patterns, access times, and user behaviour to spot deviations indicative of unauthorized entry. For example, ML algorithms trained on user behaviour profiles can identify anomalies such as unusual login locations or times, triggering alerts for potential security breaches (Géron, 2019). These models learn from legitimate patterns and dynamically update as user behaviour evolves, making them more effective than static rule-based systems.
- DDoS (Distributed Denial of Service) attacks are another pressing concern for cloud service providers. These attacks flood servers with a massive volume of requests, disrupting services and leading to potential downtime (Chukwunweike JN et al...2024). AI-driven anomaly detection systems, especially those utilizing deep learning models like LSTMs, can analyse network traffic in real time and identify patterns consistent with DDoS activities. Unlike traditional detection methods that may struggle with high-volume data or require manual threshold settings, AI techniques can recognize subtle cues leading up to an attack and flag them before the system becomes overwhelmed (Zhou et al., 2018).
- Insider threats pose a unique challenge because they originate from users with legitimate access. Detecting such threats requires systems capable of distinguishing between regular activities and malicious actions conducted under authorized credentials. Hybrid models combining unsupervised and supervised learning can monitor insider activity for deviations from established behavioural baselines without prior knowledge of specific threat signatures. This capability helps identify risky behaviour, such as unauthorized data transfers or excessive access to sensitive files, which might signal an insider threat (Rajasegarar et al., 2010).
- AI-driven anomaly detection not only enhances security by identifying these types of threats but also provides contextual insights that aid in response efforts. For instance, RL models can be programmed to recommend

actions or even autonomously respond to threats by blocking suspicious IP addresses or revoking user privileges temporarily. This adaptability supports cloud environments that require swift and intelligent responses to potential breaches. Beyond detecting and mitigating security threats, AI-driven anomaly detection can play a crucial role in other areas such as performance monitoring and resource management. By leveraging these capabilities, cloud systems can improve operational efficiency and optimize resource allocation, ensuring a well-rounded approach to cloud management.

4.2. Performance Monitoring and Optimization

4.2.1. Resource Exhaustion Detection

Anomaly detection is essential for identifying resource exhaustion, which happens when computational or memory resources become strained due to high demand or suboptimal allocation. AI-driven methods, such as unsupervised learning and deep learning, allow continuous analysis of system metrics to uncover patterns signalling impending resource shortages. Clustering algorithms, for example, can detect atypical usage trends that may lead to resource depletion (Golab et al., 2009). These insights enable proactive measures, such as reallocating resources, to prevent service disruptions during periods of peak usage.

4.2.2. Service Degradation Identification

Service degradation can severely affect user experience and occurs due to factors like network congestion, software bugs, or hardware malfunctions. Advanced AI techniques, such as autoencoders and RNN, excel at spotting subtle shifts in system performance that indicate degradation. Autoencoders, trained on historical data, compare expected and real-time performance metrics and highlight discrepancies that may signal declining service quality (Hinton & Salakhutdinov, 2006). This early warning allows IT teams to investigate and fix issues before they escalate, maintaining consistent service standards.

4.2.3. Abnormal Traffic Pattern Analysis

Monitoring for abnormal traffic patterns is crucial in cloud environments, as sudden changes can suggest both performance issues and security threats. For example, unexpected traffic spikes may result from increased user activity or malicious actions like DDoS attacks. Long Short-Term Memory (LSTM) networks are particularly adept at analysing time-series data to differentiate between routine and irregular traffic fluctuations (Hochreiter & Schmidhuber, 1997). This capability ensures that cloud providers can respond swiftly to genuine threats while avoiding unnecessary resource over-provisioning during typical traffic surges.

4.2.4. Insights for Performance Optimization

AI-based anomaly detection is not just about identifying problems but also contributing to performance optimization. Analysing the data from detected anomalies supports predictive maintenance, allowing infrastructure components to be serviced or replaced proactively (Chukwunweike JN et al., 2024). This helps minimize the risk of unplanned downtime and supports seamless operations, translating to a better user experience. With robust performance monitoring covered, the focus shifts to optimizing resource management. The next section explores how AI can enhance resource allocation strategies to improve cloud system efficiency.

4.3. AI in Resource Management and Load Balancing

AI plays a crucial role in resource management and load balancing within cloud environments, enabling optimal performance, efficiency, and resource utilization. As cloud infrastructures grow more complex, AI techniques are increasingly employed to automate the allocation of computing resources and manage fluctuating workloads in real time.

4.3.1. Dynamic Load Balancing

Dynamic load balancing is a critical process in cloud computing that ensures workloads are evenly distributed across available servers to prevent overloading any single resource. Traditional load balancing techniques often rely on static algorithms that use predefined rules and thresholds. However, as cloud environments become more dynamic, with varying traffic and resource demands, AI-driven methods offer a more adaptive approach.

ML algorithms, particularly RL, have shown promising results in optimizing load balancing. In RL, agents are trained to interact with the cloud system by making decisions that maximize efficiency based on real-time feedback (Chukwunweike JN et al., 2024). By learning from continuous interactions with the environment, these agents can

intelligently allocate tasks to servers, adjusting the load distribution based on the current state of the system (Konda & Tsitsiklis, 2003). Over time, RL models optimize decision-making, ensuring that workloads are balanced dynamically, which enhances cloud system performance and minimizes delays.

Deep learning models, such as convolutional neural networks (CNNs), can also be utilized to analyse large-scale cloud resource usage data and identify patterns that inform load balancing decisions. For example, CNNs can be trained on traffic patterns to predict future load distributions, allowing the system to pre-emptively allocate resources before traffic spikes occur, minimizing the chances of server overload.

4.3.2. Resource Allocation Optimization

Effective resource allocation ensures that cloud systems deliver high performance without unnecessary overhead. AI techniques, particularly supervised learning and optimization algorithms, are widely used to allocate resources dynamically based on demand and workload requirements. In supervised learning models, cloud administrators provide labelled data, such as historical usage patterns, which is used to train AI systems. These systems can predict resource demand and allocate resources accordingly, minimizing underutilization or overutilization (Yang et al., 2011). For instance, a supervised learning model could predict peak usage times for specific cloud services and automatically scale resources up or down, ensuring efficient resource usage.

Additionally, genetic algorithms (GAs), a form of optimization technique, can be applied to resource management. GAs simulates the process of natural evolution to find the optimal configuration of resources that minimizes cost and maximizes efficiency. In cloud computing, GAs can optimize the placement of virtual machines (VMs) or containers, ensuring that resources are allocated in a way that reduces latency and maximizes throughput (Deb, 2001).

4.3.3. Resource Management in Multi-Cloud and Hybrid Environments

In multi-cloud or hybrid environments, where resources are spread across different providers, AI can help manage and allocate resources seamlessly between different platforms. By analysing metrics from multiple cloud providers, AI systems can dynamically shift workloads between clouds based on resource availability, cost-efficiency, and performance metrics. This ensures that workloads are distributed optimally, preventing performance bottlenecks and minimizing costs.

AI-driven optimization algorithms can analyse the benefits and trade-offs of using different cloud providers and make decisions about where to place workloads based on service level agreements (SLAs), pricing models, and resource availability. This multi-cloud resource management allows businesses to leverage the strengths of various cloud providers, ensuring the best possible performance and cost efficiency. While dynamic load balancing and resource allocation optimize cloud performance, AI-driven anomaly detection plays an essential role in maintaining the reliability of these systems by identifying deviations that could signal performance issues or failures.

5. Challenges in implementing ai-driven anomaly detection in cloud computing

5.1. Data Quality and Availability

Data quality and availability are fundamental challenges in training AI models for anomaly detection in cloud environments. For AI systems to effectively detect anomalies, they require high-quality, consistent, and abundant data. However, in cloud computing, the vastness and complexity of the data involved can present significant obstacles in ensuring its quality and availability.

5.1.1. Data Quality Challenges

One of the primary concerns in cloud environments is the heterogeneity of data. Cloud systems aggregate data from multiple sources, including virtual machines, network traffic, storage systems, and user activity logs. This data can vary greatly in format, scale, and relevance, making it difficult to create a unified dataset for training AI models. For instance, data from different cloud providers or services may be structured in different ways, requiring preprocessing to standardize it for analysis. Poorly structured or incomplete data can severely hinder the performance of AI models, leading to inaccurate anomaly detection (Parker et al., 2016).

Additionally, noisy data, which contains irrelevant or incorrect information, can degrade the effectiveness of ML algorithms. Anomalies detected in noisy data may be false positives, which can lead to unnecessary system alerts or actions. For instance, if network traffic logs include redundant or extraneous entries, an anomaly detection model might

incorrectly classify a normal fluctuation as an anomaly (Marr, 2021). Filtering and cleaning data before it is used to train AI models is thus critical to achieving accurate results.

5.1.2. Data Availability Issues

The availability of data for training AI models is another significant challenge in cloud environments. Cloud service providers often store data in distributed locations across different geographic regions, creating barriers to data collection and analysis. These distributed data sources can be difficult to access in real time, leading to potential gaps in the dataset. For anomaly detection models that rely on up-to-date data, any delay in data availability can reduce the system's ability to identify anomalies promptly.

Moreover, in some cloud settings, especially multi-cloud or hybrid environments, data fragmentation can occur. In such cases, data is distributed across multiple cloud providers or on-premises systems, creating complexities in data aggregation. This makes it harder to obtain a complete view of system operations and introduces delays in the detection of anomalies that span multiple platforms. AI models, therefore, require robust data integration mechanisms to efficiently combine data from various sources (Cheng et al., 2017).

5.1.3. Impact of Insufficient or Biased Data

The effectiveness of AI-driven anomaly detection heavily depends on the data used to train the models. Insufficient data—such as a lack of historical data or data from edge cases—can result in models that fail to recognize critical anomalies. Furthermore, biased data can exacerbate this issue by training models that are not generalizable across all scenarios. For example, if an AI model is trained primarily on data from low-traffic periods, it may not effectively detect anomalies that arise during peak periods of usage (Fung et al., 2020). The lack of diverse data limits the model's ability to adapt to the dynamic nature of cloud environments, where conditions can change rapidly. Given the critical role of data in training AI models, any challenges related to data quality and availability can severely impact the effectiveness of anomaly detection systems. The next section will explore the specific consequences of poor data and how these limitations can be mitigated.

5.2. Scalability and Adaptability

One of the main challenges in applying AI-driven anomaly detection to large-scale cloud infrastructures is scalability. Cloud environments are inherently dynamic and can involve thousands or even millions of devices, users, and interactions. The ability of AI models to scale effectively while maintaining high performance becomes increasingly difficult as the volume of data and the number of cloud resources grow.

5.2.1. Scalability Issues

As cloud infrastructures expand, the amount of data generated and the complexity of system interactions increase. Traditional anomaly detection models often struggle to keep up with this growth. For example, training AI models on large datasets can lead to high computational costs and long processing times, making real-time anomaly detection difficult. Data storage and management also become more challenging, as vast amounts of unstructured data from multiple sources must be processed, stored, and analysed efficiently.

Moreover, as cloud environments are often heterogeneous, consisting of different service models (IaaS, PaaS, SaaS) and cloud providers, AI models may need to process data from various sources in parallel. The distributed nature of cloud resources requires anomaly detection systems to operate across multiple nodes and platforms simultaneously, which can strain computational resources and introduce latency.

5.2.2. Adaptability Concerns

Cloud environments are highly dynamic, with workloads constantly changing due to varying user demands, system failures, and service updates. AI-driven anomaly detection models must be adaptable to these changes, but traditional models are typically static and fail to account for the evolving nature of cloud infrastructures. This can result in models that do not generalize well to new scenarios, making it difficult to detect new or unseen anomalies.

For example, if the system experiences a sudden shift in traffic patterns or a new type of service deployment occurs, an AI model trained on previous data may fail to detect anomalies associated with these changes. To address scalability and adaptability challenges, several methods and strategies are being developed. The next section will explore techniques for enhancing scalability in AI-driven anomaly detection systems.

5.3. Computational and Resource Constraints

AI-based anomaly detection models, particularly those leveraging deep learning and RL, require substantial computational resources to operate effectively in cloud environments. These models are typically data-intensive and demand significant processing power, memory, and storage, posing challenges for their deployment at scale.

5.3.1. High Computational Costs

One of the primary issues with AI-driven models in cloud environments is the computational cost. Deep learning models, for example, require powerful GPUs or TPUs to train and infer from vast datasets. As the size of cloud infrastructures grows, the volume of data to be processed also increases, leading to longer training times and higher resource consumption. For instance, training a neural network on a large cloud system dataset can take hours or even days, depending on the model complexity and data volume. These prolonged training times increase the overall cost of computation, as cloud providers often charge based on the resources consumed (e.g., GPU hours, storage) (Zhang et al., 2020).

5.3.2. Resource Limitations

Cloud environments, despite offering flexible resource allocation, still face resource limitations when it comes to AI-based anomaly detection. While virtual machines (VMs) or containers can be provisioned for running AI models, the demand for computational resources may exceed what is available, especially during peak usage periods. This can result in resource contention, where multiple tasks or services compete for limited resources, slowing down anomaly detection processes. Additionally, running complex models on multiple nodes across distributed systems increases the overhead, leading to challenges in maintaining efficiency while avoiding latency (Chen et al., 2020).

5.3.3. Storage Requirements

Another significant resource constraint is the need for large-scale storage to store historical data and trained models. Cloud systems need extensive storage capabilities to manage datasets used for training and testing models. Managing this data in real time for anomaly detection requires highly efficient data storage and retrieval systems, which can become a bottleneck if not properly optimized (Sheng et al., 2021). While computational costs and resource constraints present challenges, several optimization techniques can help address these issues. The next section will explore strategies for optimizing AI-driven models to reduce these resource demands.

6. Optimizing AI-based anomaly detection models for cloud environments

6.1. Model Efficiency and Optimization

The efficiency of AI models plays a crucial role in their effectiveness within cloud environments. AI-based anomaly detection systems must process large volumes of data in real time, which can be computationally intensive and resource-heavy. As cloud infrastructures grow in size and complexity, optimizing these models becomes imperative for ensuring fast, efficient anomaly detection without overburdening the system's resources. Several optimization techniques have been developed to improve the processing speed, reduce resource consumption, and maintain the accuracy of AI models used for anomaly detection.

6.1.1. Model Pruning

One of the primary optimization techniques for AI models is model pruning. Model pruning involves the removal of unnecessary or redundant parameters in neural networks, reducing the number of weights and, consequently, the complexity of the model. This technique focuses on identifying and eliminating neurons or connections that do not contribute significantly to the model's performance. By pruning the model, it becomes more efficient in terms of computational resources and can operate faster, which is crucial for real-time anomaly detection (Han et al., 2015).

Pruning is particularly effective in cloud environments, where AI models need to handle large, high-dimensional datasets. By reducing the number of parameters, the model's memory footprint is decreased, allowing it to run on devices or cloud nodes with limited resources, without sacrificing performance. This also translates to quicker inference times, which is critical for detecting anomalies promptly before they escalate into more significant issues.

6.1.2. Quantization

Another widely used optimization technique is quantization. Quantization reduces the precision of the numerical values used in the model's parameters, typically transforming them from 32-bit floating-point numbers to lower-bit

representations, such as 8-bit integers. This leads to a reduction in both memory usage and computation time. While quantization may slightly reduce the model's accuracy, careful tuning and post-training quantization techniques can minimize this loss and still maintain the model's effectiveness (Banner et al., 2019).

In cloud environments, where multiple users or services may be running concurrently, quantization helps ensure that AI models can be deployed more efficiently without requiring excessive computational power. It also allows for faster model inference, which is important in dynamic environments where quick responses are needed, such as detecting abnormal network traffic or unauthorized access in real-time.

6.1.3. Transfer Learning

Transfer learning is another powerful optimization technique that can be used to improve the efficiency of AI models for anomaly detection in cloud environments. Transfer learning involves leveraging a pre-trained model on a similar problem and fine-tuning it on the specific dataset of interest. This method reduces the computational cost of training a new model from scratch by using knowledge from an already trained model (Pan & Yang, 2010). In cloud environments, this approach is particularly beneficial because it minimizes the need for large-scale data processing and extensive training resources.

By fine-tuning a pre-trained model, cloud-based anomaly detection systems can be adapted to detect specific anomalies related to a particular service or system without the need for training an entirely new model. This can save significant time and computational resources while ensuring that the system remains effective in identifying potential threats and performance issues.

6.1.4. Optimizing Inference

Beyond model-specific optimizations, inference optimization also plays a key role in ensuring the overall efficiency of AI-driven anomaly detection systems. Cloud environments often have varying levels of available resources, and optimizing the inference process ensures that the model can operate efficiently under different conditions. Techniques like batch processing, where data is processed in groups rather than individually, and model distillation, which creates smaller, more efficient versions of large models, can significantly improve inference speed while maintaining accuracy (Hinton et al., 2015).

Furthermore, specialized hardware accelerators, such as FPGAs or TPUs, can be utilized to enhance inference speed and reduce energy consumption. These hardware accelerators are optimized for running AI models, enabling faster anomaly detection and improved scalability, even in highly complex cloud environments.

6.1.5. Efficient Data Handling

Efficient data handling is also an essential part of optimizing AI-driven anomaly detection in cloud environments. Data preprocessing and feature extraction are critical to ensuring that the AI model processes only the most relevant and useful information. This can reduce the amount of data that needs to be handled and processed by the model, further lowering computational costs and speeding up inference.

Cloud environments often deal with heterogeneous data, meaning that multiple data sources, each with different formats, must be processed. Optimizing the way data is ingested, cleaned, and transformed before being fed into an AI model is crucial for reducing the overall resource load. For example, data aggregation techniques can be used to summarize or compress data, providing the model with essential information while reducing the total data volume. By employing these model efficiency and optimization techniques, cloud environments can leverage AI-driven anomaly detection systems that are not only effective but also efficient and resource-conscious. As cloud infrastructures grow in complexity, optimizing AI models becomes even more critical for ensuring real-time anomaly detection without compromising performance. The next section will discuss how these optimizations can be integrated into practical applications of AI-based anomaly detection systems.

6.2 Federated Learning and Edge Computing

Federated learning and edge computing are innovative technologies that can enhance the efficiency and scalability of AI-driven anomaly detection systems, particularly in cloud environments. By decentralizing processing tasks, these approaches offer significant benefits in handling large datasets and reducing the burden on centralized cloud infrastructure.

6.1.6. Federated Learning

Federated learning is a ML paradigm where multiple devices or nodes collaboratively train a model while keeping their data local. In the context of anomaly detection, federated learning allows edge devices, such as IoT sensors or local servers, to perform local training without needing to send raw data to the cloud. This approach improves efficiency by reducing data transfer and maintaining privacy, as only model updates (rather than the actual data) are shared between devices and the central server (McMahan et al., 2017).

In cloud-based anomaly detection, federated learning can be particularly useful for detecting distributed anomalies across diverse endpoints. This decentralized approach allows the system to adapt to different local data distributions, leading to more personalized and accurate anomaly detection models. Moreover, by processing data closer to the source, federated learning reduces latency and increases the overall responsiveness of the system, which is critical in real-time anomaly detection scenarios.

6.1.7. Edge Computing

Edge computing further complements federated learning by decentralizing the computational tasks associated with anomaly detection. Edge computing refers to processing data closer to the source—on edge devices, rather than relying on centralized cloud infrastructure. By moving processing power to the edge, anomaly detection tasks are executed locally, enabling faster detection of issues such as network intrusions or service disruptions. This is particularly beneficial in environments where real-time responses are crucial, as it reduces the need for data transmission and can mitigate delays caused by cloud processing.

Together, federated learning and edge computing allow AI-driven anomaly detection systems to scale efficiently, as the workload is distributed across multiple devices. This enhances the system's overall performance and robustness, especially in large and dynamic cloud infrastructures. While decentralizing processing through federated learning and edge computing enhances efficiency, it is crucial to evaluate the performance of these systems to ensure their effectiveness. The next section will focus on the evaluation metrics used to assess AI-driven anomaly detection models.

6.2. Evaluation Metrics and Model Assessment

Assessing the performance of AI-driven anomaly detection systems is essential to ensure their reliability and effectiveness in cloud environments. Several key evaluation metrics help determine how well these models perform in identifying anomalies while minimizing false positives and negatives.

6.2.1. Precision and Recall

Precision and recall are two fundamental metrics in evaluating anomaly detection systems. Precision refers to the percentage of correctly identified anomalies out of all the instances flagged as anomalies, while recall measures the percentage of actual anomalies correctly identified by the system. These two metrics help balance the trade-off between false positives (incorrectly flagged normal data) and false negatives (missed anomalies). In cloud environments, where the cost of both types of errors can be significant, it is important to find a balance that minimizes both false positives and false negatives (Chandola et al., 2009).

6.2.2. F1 Score

The F1 score combines precision and recall into a single metric, providing a harmonic mean of the two. It is particularly useful in scenarios where there is an uneven class distribution, such as when anomalies are rare events in a cloud system. A high F1 score indicates a good balance between precision and recall, making it an important metric for anomaly detection systems (Saito & Rehmsmeier, 2015).

6.2.3. ROC Curve

The Receiver Operating Characteristic (ROC) curve is another important tool for evaluating anomaly detection models. It plots the true positive rate (sensitivity) against the false positive rate (1-specificity), providing a visual representation of a model's performance across different thresholds. The area under the ROC curve (AUC) serves as a summary statistic for the model's overall performance. A higher AUC value indicates better discrimination between normal and anomalous behaviour (Fawcett, 2006).

6.2.4. Other Metrics

In addition to these core metrics, other evaluation measures, such as AUC-PR (Area Under the Precision-Recall Curve) and confusion matrices, may also be used to gain a more comprehensive understanding of the model's performance in detecting anomalies, particularly in imbalanced datasets common in cloud environments. These evaluation metrics are crucial for assessing the effectiveness of anomaly detection systems. Looking forward, the next section will explore future trends and research directions in AI-driven anomaly detection for cloud environments.

7. Future trends and research directions in AI-driven anomaly detection for cloud computing

7.1. Advancements in AI Techniques

The field of artificial intelligence (AI) is rapidly evolving, introducing new techniques that hold promise for enhancing anomaly detection in cloud environments. Among these, self-supervised learning and explainable AI (XAI) are particularly noteworthy due to their potential to improve the accuracy, transparency, and interpretability of AI-driven anomaly detection systems.

7.1.1. Self-Supervised Learning

Self-supervised learning is an emerging ML paradigm that enables models to learn useful representations from unlabelled data by generating pseudo-labels through intrinsic patterns within the data itself. This technique is particularly beneficial in cloud environments where acquiring labelled data can be challenging due to the sheer volume and complexity of cloud-based activities. Self-supervised learning models can be trained on massive, unlabelled datasets and still learn to detect anomalies effectively. For instance, a self-supervised model could identify unusual traffic patterns or resource utilization trends in cloud systems by leveraging temporal or spatial correlations in the data (Berthelot et al., 2019). By utilizing vast amounts of unlabelled data, this technique can address data scarcity issues, making it a valuable tool in the cloud environment where labelled datasets are often limited or expensive to curate.

7.1.2. XAI

XAI aims to make AI models more interpretable by providing human-readable explanations for their predictions and decisions. This is especially crucial in anomaly detection, where the system must not only detect anomalies but also provide insight into why certain behaviours are considered anomalous. Cloud systems often involve complex and opaque processes, making it difficult for security analysts to trust black-box AI models. XAI can help bridge this gap by providing transparent reasoning behind the detection of anomalies, thus improving trust and aiding decision-making (Gilpin et al., 2018). For example, an XAI model might explain that an anomaly was detected due to a sudden surge in inbound traffic or the use of a previously unused API key, making it easier for security teams to understand and mitigate the threat. The advancements in AI techniques like self-supervised learning and XAI are part of a broader trend of increasing automation and intelligence in cloud security. As cloud environments become more complex, these innovations provide the foundation for next-generation anomaly detection systems. The integration of AI-driven anomaly detection with emerging technologies further amplifies its potential.

7.2. Integration with Emerging Technologies

AI-driven anomaly detection systems in cloud environments can be further enhanced by integrating with other emerging technologies such as blockchain and Internet of Things (IoT). These technologies offer additional layers of security, reliability, and scalability, complementing the capabilities of AI.

7.2.1. Blockchain for Data Integrity

Blockchain technology, known for its decentralized and tamper-resistant nature, can be integrated with AI-driven anomaly detection to enhance data integrity and auditability in cloud environments. Blockchain can provide a secure ledger of all actions and transactions within a cloud system, ensuring that the data fed into AI models remains unaltered and trustworthy. By using blockchain to verify data sources, AI models can detect anomalies not only in real-time operations but also by cross-checking historical data. Moreover, the immutability of blockchain can help in creating transparent, immutable logs of detected anomalies, which can be useful for compliance and regulatory purposes (Narayanan et al., 2016).

7.2.2. Internet of Things (IoT) Integration

As the Internet of Things (IoT) continues to proliferate, integrating IoT devices with cloud-based anomaly detection systems becomes increasingly important. IoT devices generate a constant stream of data, which, when analysed in real-

time, can help detect anomalies in cloud infrastructures, such as unexpected patterns of resource consumption or unauthorized access attempts. By connecting IoT sensors to AI-driven anomaly detection systems, cloud environments can be monitored for potential security breaches, system failures, or inefficiencies. Furthermore, the vast volume and variety of data generated by IoT devices necessitate the use of advanced AI techniques to process and analyse this information, making IoT integration a key consideration in future cloud security solutions (Muthusamy & Prabadevi, 2020). While the integration of AI-driven anomaly detection with emerging technologies like blockchain and IoT holds great promise, it also presents several challenges. The following section will discuss some of the key challenges and opportunities for future research in this area, including ethical and regulatory considerations.

7.3. Challenges and Opportunities for Future Research

Despite the promising advancements and integration opportunities, several challenges remain in the development and implementation of AI-driven anomaly detection systems in cloud environments. Addressing these challenges will be critical for ensuring the continued evolution of these technologies.

7.3.1. Ethical and Regulatory Challenges

One of the major challenges in deploying AI-driven anomaly detection in cloud environments is the ethical and regulatory concerns surrounding data privacy, security, and bias. AI models require vast amounts of data to be trained effectively, and in the case of cloud systems, this data can include sensitive information such as user behaviour, financial transactions, and personal data. Protecting this data while ensuring that anomaly detection systems are effective is a delicate balance. Regulations such as the General Data Protection Regulation (GDPR) in the European Union impose strict requirements on how personal data is collected, processed, and stored. Future research must explore ways to develop anomaly detection systems that comply with these regulations, ensuring data privacy and security while maintaining effectiveness (Voigt & Von dem Bussche, 2017).

Moreover, AI models are susceptible to biases that may arise from unbalanced datasets or flawed training processes. This can lead to skewed anomaly detection results, such as a higher number of false positives for certain user groups or activities. Ensuring fairness and transparency in AI models is an ongoing area of research that must be addressed to build trust in AI-driven security solutions.

7.3.2. Opportunities for Research

The integration of edge computing, 5G networks, and quantum computing presents new opportunities for future research in AI-driven anomaly detection. Edge computing can decentralize anomaly detection, reducing latency and improving real-time responses in cloud environments. Additionally, the advent of 5G networks offers increased bandwidth and lower latency, which can enhance the effectiveness of real-time anomaly detection systems, particularly in high-traffic cloud environments. Quantum computing, although still in the early stages, holds the potential to revolutionize anomaly detection by enabling faster processing of large-scale datasets, allowing AI models to detect more complex anomalies with greater precision (Arute et al., 2019).

Therefore, the future of AI-driven anomaly detection in cloud environments is rife with opportunities and challenges. As AI techniques evolve and integrate with emerging technologies, the potential for enhanced cloud security grows. However, addressing ethical, regulatory, and technical challenges will be crucial for the successful deployment of these systems. Ongoing research will be essential to overcome these hurdles, ensuring that AI-driven anomaly detection can effectively safeguard cloud environments while maintaining privacy, fairness, and transparency.

8. Conclusion

This research has explored the critical role of AI-driven anomaly detection in enhancing the security, performance, and reliability of cloud computing environments. As cloud infrastructure continues to expand in complexity and scale, the need for advanced security measures has become paramount. Traditional anomaly detection methods are no longer sufficient to address the dynamic and highly distributed nature of cloud environments. In this context, AI techniques, particularly ML, have emerged as a powerful solution, offering more efficient, scalable, and accurate detection of unusual patterns and potential security threats.

One of the key findings from this paper is that supervised, unsupervised, and semi-supervised learning methods provide a comprehensive framework for detecting both known and unknown anomalies in cloud systems. AI's ability to analyse vast amounts of real-time data from diverse sources, such as servers, network traffic, and user behaviour, allows for the identification of security breaches, system failures, and resource inefficiencies with much higher precision than

traditional methods. Specifically, deep learning, including techniques such as autoencoders and RNNs, has proven particularly effective in handling the complexity of cloud environments by learning from historical data to detect both subtle and significant deviations from normal behaviour.

Furthermore, XAI and self-supervised learning have introduced significant improvements in the interpretability and scalability of anomaly detection systems. XAI models, which provide human-readable explanations of detected anomalies, help cloud security teams understand the underlying causes of detected threats, facilitating faster response and resolution. Self-supervised learning allows for the training of models on large amounts of unlabelled data, addressing the challenge of acquiring labelled datasets in cloud environments, where data is often unstructured or too vast to label manually. These advances in AI techniques have significantly increased the robustness of anomaly detection systems, making them more adaptable to the ever-evolving nature of cloud infrastructures.

The integration of AI-driven anomaly detection with emerging technologies like blockchain and the Internet of Things (IoT) has opened up new possibilities for improving cloud security. Blockchain provides an immutable record of activities within cloud environments, ensuring data integrity and providing transparency in anomaly detection processes. On the other hand, IoT integration enhances anomaly detection by leveraging real-time data from IoT devices, enabling faster and more accurate identification of security threats or system anomalies. These technologies, when combined with AI, create a more secure and resilient cloud infrastructure, capable of responding to both known and unknown threats in real time.

Despite the promising advancements in AI-driven anomaly detection, several challenges remain. Ethical considerations, such as data privacy and fairness, must be addressed to ensure that AI models do not inadvertently violate user privacy or discriminate against certain user groups. Regulatory frameworks like the General Data Protection Regulation (GDPR) impose strict rules on the use of personal data, which AI-driven systems must adhere to in order to maintain trust and compliance. Additionally, the scalability and adaptability of AI models remain significant challenges, particularly when applied to large and dynamic cloud infrastructures. As cloud systems grow and evolve, anomaly detection systems must be able to scale accordingly without compromising performance or resource efficiency.

Looking forward, the future of AI in cloud computing is bright, with several exciting opportunities on the horizon. The advent of edge computing and 5G networks will further enhance the performance of AI-driven anomaly detection by reducing latency and enabling real-time processing of vast amounts of data at the edge of the network. Federated learning, which allows AI models to be trained across decentralized devices while keeping data local, will also address data privacy concerns and improve the scalability of anomaly detection systems. Moreover, quantum computing, though still in its infancy, holds the potential to revolutionize AI anomaly detection by enabling faster processing and more accurate anomaly detection on large datasets.

Thus, AI-driven anomaly detection is poised to become a cornerstone of cloud security in the coming years. As AI technologies continue to evolve, their integration with emerging technologies like edge computing, 5G, blockchain, and IoT will create increasingly intelligent and adaptive security systems. However, addressing the challenges of data privacy, model interpretability, and scalability will be crucial for ensuring that AI-driven anomaly detection systems are effective, reliable, and ethical. The continued research and development of these AI techniques will play a critical role in shaping the future of cloud computing, ensuring that cloud environments remain secure, resilient, and efficient in the face of evolving cyber threats.

Compliance with ethical standards

Disclosure of conflict of interest

No conflict of interest to be disclosed.

Reference

- [1] Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3), 1-58. <https://doi.org/10.1145/1541880.1541882>
- [2] Iglewicz, B., & Hoaglin, D. C. (1993). *How to detect and handle outliers*. Sage Publications.

- [3] Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., & Brandic, I. (2010). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation Computer Systems*, 25(6), 599-616. <https://doi.org/10.1016/j.future.2008.12.001>
- [4] Cai, Z., Zhang, Z., Lin, M., Zheng, Y., & Yu, Y. (2016). A data-driven approach for real-time insider threat detection. *IEEE Transactions on Dependable and Secure Computing*, 14(2), 205-216. <https://doi.org/10.1109/TDSC.2016.2565566>
- [5] Greitzer, F. L., & Frincke, D. A. (2010). Combining traditional cyber security audit data with user behavior analysis to improve insider threat detection. *Journal of Information Warfare*, 8(3), 1-15.
- [6] Mell, P., & Grance, T. (2011). The NIST definition of cloud computing. *NIST Special Publication*, 800-145. <https://doi.org/10.6028/NIST.SP.800-145>
- [7] Mirkovic, J., & Reiher, P. (2004). A taxonomy of DDoS attack and DDoS defense mechanisms. *ACM SIGCOMM Computer Communication Review*, 34(2), 39-53. <https://doi.org/10.1145/997150.997156>
- [8] Zhang, R., Liu, L., & Chen, J. (2010). Privacy-preserving cloud computing: Challenges and opportunities. *IEEE Data Engineering Bulletin*, 32(1), 27-34.
- [9] Xu, X., & Tian, Y. (2015). A comprehensive survey of clustering algorithms. *Annals of Data Science*, 2(2), 165-193. <https://doi.org/10.1007/s40745-015-0040-1>
- [10] Hodge, V. J., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2), 85-126. <https://doi.org/10.1023/B:AIRE.0000045502.10941.a9>
- [11] Xu, X., & Tian, Y. (2015). A comprehensive survey of clustering algorithms. *Annals of Data Science*, 2(2), 165-193. <https://doi.org/10.1007/s40745-015-0040-1>
- [12] Liu, F. T., Ting, K. M., & Zhou, Z.-H. (2008). Isolation forest. *Proceedings of the 2008 IEEE International Conference on Data Mining*, 413-422. <https://doi.org/10.1109/ICDM.2008.17>
- [13] Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504-507. <https://doi.org/10.1126/science.1127647>
- [14] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [15] Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press.
- [16] Joseph Nnaemeka Chukwunweike, Moshood Yussuf, Oluwatobiloba Okusi, Temitope Oluwatobi Bakare, Ayokunle J. Abisola. The role of deep learning in ensuring privacy integrity and security: Applications in AI-driven cybersecurity solutions [Internet]. Vol. 23, World Journal of Advanced Research and Reviews. GSC Online Press; 2024. p. 1778-90. Available from: <https://dx.doi.org/10.30574/wjarr.2024.23.2.2550>
- [17] Zhou, Z.-H., Wu, J., & Tang, W. (2002). Ensembling neural networks: Many could be better than all. *Artificial Intelligence*, 137(1-2), 239-263. [https://doi.org/10.1016/S0004-3702\(02\)00190-X](https://doi.org/10.1016/S0004-3702(02)00190-X)
- [18] Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems* (2nd ed.). O'Reilly Media.
- [19] Zhou, R., Li, Z., & Chen, T. (2018). Anomaly detection with LSTM networks in cloud environments. *Journal of Cloud Computing*, 7(1), 1-15. <https://doi.org/10.1186/s13677-018-0123-4>
- [20] Rajasegarar, S., Leckie, C., & Palaniswami, M. (2010). Anomaly detection in wireless sensor networks. *IEEE Journal on Selected Areas in Communications*, 28(7), 1424-1435. <https://doi.org/10.1109/JSAC.2010.100701>
- [21] Golab, L., Demaine, E. D., & McGregor, A. (2009). *Stream data processing and resource allocation in cloud environments*. *IEEE Transactions on Computers*, 58(9), 1207-1216. <https://doi.org/10.1109/TC.2009.63>
- [22] Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504-507. <https://doi.org/10.1126/science.1127647>
- [23] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [24] Konda, V. R., & Tsitsiklis, J. N. (2003). *Actor-Critic Algorithms*. *Advances in Neural Information Processing Systems*, 12, 1008-1014.

- [25] Chukwunweike JN, Adewale AA, Osamuyi O 2024. Advanced modelling and recurrent analysis in network security: Scrutiny of data and fault resolution. DOI: [10.30574/wjarr.2024.23.2.2582](https://doi.org/10.30574/wjarr.2024.23.2.2582)
- [26] Yang, Q., Liu, Y., & Tang, X. (2011). Optimizing resource allocation in cloud computing environments. *Journal of Cloud Computing: Advances, Systems and Applications*, 1(1), 1-17. <https://doi.org/10.1186/2192-113X-1-1>
- [27] Banner, R., Korman, G., & Shkolnik, M. (2019). *Post-training quantization for deep neural networks*. *IEEE Transactions on Neural Networks and Learning Systems*, 30(5), 1304-1316. <https://doi.org/10.1109/TNNLS.2018.2833243>
- [28] Han, S., Mao, H., & Dally, W. J. (2015). *Deep compression: Compressing deep neural networks with pruning, trained quantization, and Huffman coding*. *Proceedings of the International Conference on Learning Representations (ICLR)*. <https://openreview.net/forum?id=BJt3qS5gg>
- [29] Hinton, G., Vinyals, O., & Dean, J. (2015). *Distilling the knowledge in a neural network*. *Proceedings of the NIPS Deep Learning and Representation Learning Workshop*. <https://arxiv.org/abs/1503.02531>
- [30] Chukwunweike JN, Praise A, Bashirat BA, 2024. Harnessing Machine Learning for Cybersecurity: How Convolutional Neural Networks are Revolutionizing Threat Detection and Data Privacy. <https://doi.org/10.55248/gengpi.5.0824.2402>.
- [31] Pan, S. J., & Yang, Q. (2010). *A survey on transfer learning*. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345-1359. <https://doi.org/10.1109/TKDE.2009.191>
- [32] Fawcett, T. (2006). *An introduction to ROC analysis*. *Pattern Recognition Letters*, 27(8), 861-874. <https://doi.org/10.1016/j.patrec.2005.10.010>
- [33] McMahan, H. B., Moore, E., & Ramage, D. (2017). *Communication-efficient learning of deep networks from decentralized data*. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS 2017)*. <https://arxiv.org/abs/1602.05629>
- [34] Saito, T., & Rehmsmeier, M. (2015). *The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets*. *PLOS ONE*, 10(3), e0118432. <https://doi.org/10.1371/journal.pone.0118432>
- [35] Arute, F., et al. (2019). *Quantum supremacy using a programmable superconducting processor*. *Nature*, 574(7779), 505-510. <https://doi.org/10.1038/s41586-019-1666-5>
- [36] Chukwunweike JN, Praise A, Osamuyi O, Akinsuyi S and Akinsuyi O, 2024. AI and Deep Cycle Prediction: Enhancing Cybersecurity while Safeguarding Data Privacy and Information Integrity. <https://doi.org/10.55248/gengpi.5.0824.2403>
- [37] Berthelot, D., et al. (2019). *Self-supervised learning of pretext-invariant representations*. *Proceedings of the 36th International Conference on Machine Learning (ICML 2019)*. <https://arxiv.org/abs/1905.05594>
- [38] Gilpin, L. H., et al. (2018). *Explaining explanations: An overview of interpretability of machine learning*. *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. <https://doi.org/10.1145/3173574.3174113>
- [39] Muthusamy, R., & Prabadevi, P. (2020). *AI-powered anomaly detection in IoT systems*. *IEEE Access*, 8, 213334-213343. <https://doi.org/10.1109/ACCESS.2020.3031237>
- [40] Joseph Nnaemeka Chukwunweike and Opeyemi Aro. Implementing agile management practices in the era of digital transformation [Internet]. Vol. 24, World Journal of Advanced Research and Reviews. GSC Online Press; 2024. Available from: DOI: [10.30574/wjarr.2024.24.1.3253](https://doi.org/10.30574/wjarr.2024.24.1.3253)
- [41] Narayanan, A., Bonneau, J., Felten, E., Miller, A., & Shacham, H. (2016). *Bitcoin and cryptocurrency technologies*. Princeton University Press.
- [42] Voigt, P., & Von dem Bussche, A. (2017). *The EU General Data Protection Regulation (GDPR)*. Springer Vieweg, Berlin, Heidelberg. <https://doi.org/10.1007/978-3-662-56797-3>
- [43] *Multi-Objective Optimization Using Evolutionary Algorithms*. Wiley. <https://doi.org/10.1002/9781119191402>
- [44] Cheng, X., Zhang, M., & Zhao, X. (2017). *A survey of multi-cloud data management techniques and challenges*. *Future Generation Computer Systems*, 68, 94-107. <https://doi.org/10.1016/j.future.2016.08.022>
- [45] Fung, W., Yao, Y., & Lin, W. (2020). *Bias and fairness in machine learning: A review of the landscape*. *Journal of Cloud Computing: Advances, Systems, and Applications*, 9(1), 1-15. <https://doi.org/10.1186/s13677-020-00199-5>

- [46] Marr, B. (2021). *Data-driven: Creating a data culture*. Wiley.
- [47] Parker, G., Dawson, T., & Young, M. (2016). *The challenges of big data and machine learning in cloud environments*. *International Journal of Cloud Computing and Services Science*, 5(3), 149-162. <https://doi.org/10.1007/s11761-016-0201-2>
- [48] Chen, G., Sun, Z., & Wu, Y. (2020). *Resource allocation for AI-based cloud computing systems*. *Journal of Cloud Computing: Advances, Systems, and Applications*, 9(1), 1-10. <https://doi.org/10.1186/s13677-020-00240-w>
- [49] Sheng, Q. Z., Li, C., & Zhang, X. (2021). *Efficient storage management for AI workloads in cloud environments*. *International Journal of Computer Applications*, 43(3), 142-155. <https://doi.org/10.1080/01446193.2021.1877587>
- [50] Adetunji, A. S., Afolayan, A., Olola, T., Fonkem, B., & Odunayo, R. (2023). Enhancing STEM Education through Culturally Relevant Engineering Design: A Mixed-Methods Approach to Improving Student Retention and Engagement. <https://doi.org/10.5281/zenodo.14018509>
- [51] Zhang, Jauro F, Chiroma H, Gital AY, Almutairi M, Shafi'i MA, Abawajy JH. Deep learning architectures in emerging cloud computing architectures: Recent development, challenges and next research trend. *Applied Soft Computing*. 2020 Nov 1;96:106582. <https://doi.org/10.1016/j.asoc.2020.106582>