# Observability Driven Incident Management for Cloud-native Application Reliability

## Anila Gogineni

Independent Researcher, USA
anila.ssn@gmail.com

**Abstract**

**Cloud-native indeed tends to represent a new generation of software applications based on such principles as microservices, containerization, and real-time orchestration to achieve scalability, flexibility, and redundancy. However, these distributed systems add up considerable operational overheads due to system interdependence, workload fluctuations and frequently changing system status. These in turn often cause problems in identifying, diagnosing and solving these incidents, which in turn impacts the reliability of the service being offered. Through the lens of observability, we see monitoring as a way to solve complexities of cloud-native systems through the ability to get insights into their workings. The major difference between monitoring and observability is that while monitoring few metrics and thresholds, observing trains the capability to have visibility into the inner workings of systems on logs, metrics, and traces. It makes it easier to prevent issues before they arise, find out why they have occurred and close incidents quickly. The goal of this report is to analyze the part that observability plays in improving the reliability of cloud-native applications by using Incident Management as a prism. It puts forward a scalable model that folds observability into the experience of defining incidents and correlating them and applying remediation autonomously. Using key principles, existing tools, and practical implementations, the report proves that observability is revolutionizing MTTR and System Uptime.**

**Keywords: Cloud-Native Applications, Microservices Architecture, Observability, Root Cause Analysis, Service Monitoring, Kubernetes, Grafana, Cloud Monitoring, Data Pipelines, Cloud Infrastructure.**

## I. INTRODUCTION

Applications which are natively built for the cloud are of a new generation of software solutions because they use microservices architecture, containers and dynamic orchestration. These applications exist as distributed systems where the parts of the overall application can work automatically but will send and receive messages through the network. Because of the aforementioned strengths, this modular approach brings in complications like dependency among services, dynamic utilization of resources as well as constant surveillance of multiple components. Assim, observability, that is the capacity to ascertain the internal state of a system by measuring the output signals it selective provides, has become a starting point for guaranteeing cloud-native reliable systems. Given that observability comprises logs, metrics, and traces, it is easy for teams to identify anomalies and work on the root course to solve issues. This is especially important in today's cloud native environment where terminating an incident often does not solve the problem. Traditional approaches that depend on static key performance indicators and interventions cannot cope with currently complex and volatile systems causing long outages and higher costs. This report looks at the importance of observability-driven management of incidents in improving the resilience of cloud-native

apps. It reviews fundamental issues associated with these systems, explains what observability is and how it differs from other concepts, and presents the framework for adopting observability across the incident management life cycle.

## II. THE CHALLENGES OF INCIDENT MANAGEMENT IN CLOUD-NATIVE SYSTEMS

### Distributed Architectures and Microservices Complexity

Cloud-native applications are built using microservices architecture that splits the application into a number of different modules each of which delivers some part of the application functionality. Such an approach enables factorization of independent components and short firm cycles of deployment. Nevertheless, it also creates many operational problems. These services operate within their own container, always communicating with other services over networks as well as external services such as databases or third party APIs [1].

The identified issues are further exacerbated by the emerging characteristics of cloud-native applications. Being services, they are able to be sized up or down and containers can be started or destroyed in a matter of seconds.

### Common Types of Incidents

There are several types of problems that can affect cloud-native applications, these include; delay, contention for resources and network failures.

*A. Latency Issues*

Solutions with certain levels of response time issues or component inability to provide necessary performance thresholds might provide poor experience for the user due to overloaded services or main flows that interact with components inadequately.

*B. Resource Contention*

When the availability of hardware resources such as compute power, memory, disk or databases are restricted it causes definite constraints particularly while there is higher traffic. [2].

*C. Network Failures*

When there are many microservices interacting with each other over a network, then this produced gap can greatly affect the entire application. These problems are worsened by intermittent faults which make it difficult both to diagnose and correct the situation.
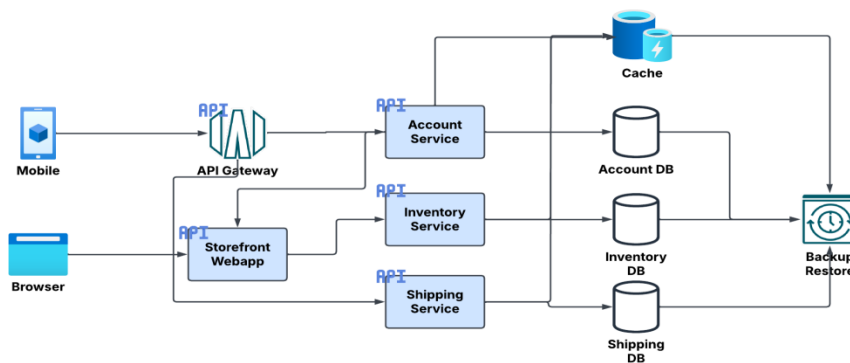


**FIG 1: The Interdependencies of Microservices, Communication Layers, And Potential Fault Zones**

**Challenges in Root-Cause Analysis**

It is, therefore, relatively difficult to find out the cause of certain incidents especially in the cloud-native systems. This means one would think of a service, and a single failure might be a domino effect in different services, so one could hardly tell where the problem started. Moreover, the absence of a single observability approach frequently leads to dispersed data. Some services might generate log, metrics, and traces in an unrelated manner that needed to be mobilized manually in order to gain an understanding of the incident.

Traditional spread sheet reviews that rely on variance analysis and other ad hoc checkpoints are inadequate in these circumstances. For instance, basic performance tools can discover problems, for example, CPU intensity or memory use yet will overlook the underlying problem which leads to protracted downtime [13].

**Limitations of Existing Incident Management Practices**

The initial approaches to handle the incidents are not suitable in Cloud-native context as it is based on the Centralized and Vertically Integrated Structure. These strategies generally require actions based on set alarms and consequent preventive procedures. The above solutions do not suffice for real-time anomaly detection and preventive incident handling in dynamic distributed systems [3]. Also, the lack of automated tools to correlate metrics, traces, and logs results in longer resolution times and a high mean-time-to-resolution (MTTR).

## III. OBSERVABILITY CONCEPTS AND THEIR ROLE IN RELIABILITY

**Defining Observability in Software Systems**

In the context of software designs, "observability" means how much one can deduce about what is going on inside a system based solely on what is coming out of it, a concept originally borrowed from control theory, and is an essential mechanism for analysis, diagnostics, and optimization of cloud-native applications. Compared to purely monitoring, observability means not only noticing symptoms of problems but offering the ability to identify the root cause by providing deep insight into system behavior. Of particular importance as a reliability enabler is its ability to support reliability in complex environments where networks are dynamic and distributed.

**Core Pillars of Observability**

The foundation of observability rests on three core pillars: metrics, logs, and traces. These elements work together to provide a complete picture of system performance and health:

A. *Metrics:*

Numerical data represented over time summarizing the numerical values provided by each of the system components as regards to their work load and efficiency according to the likes of CPU use, amount of memory used or number of requests per time interval [4]. Measures allow an evaluation of trends and regularities, making it easier for teams to consider 'abnormal' behavior.

B. *Logs:*

Contemporary asynchronous log records that contain textual information on specific occurrences in the system, failures, or events. Records hold a significant value for investigative analysis and a supportive contextual investigation environment.

*C. Traces:*

Distributed traces represent the flow of requests in between services, including latency and other characteristics of performance. Only thanks to traces, it is possible to define bottlenecks and analyze the dependencies between the microservices.
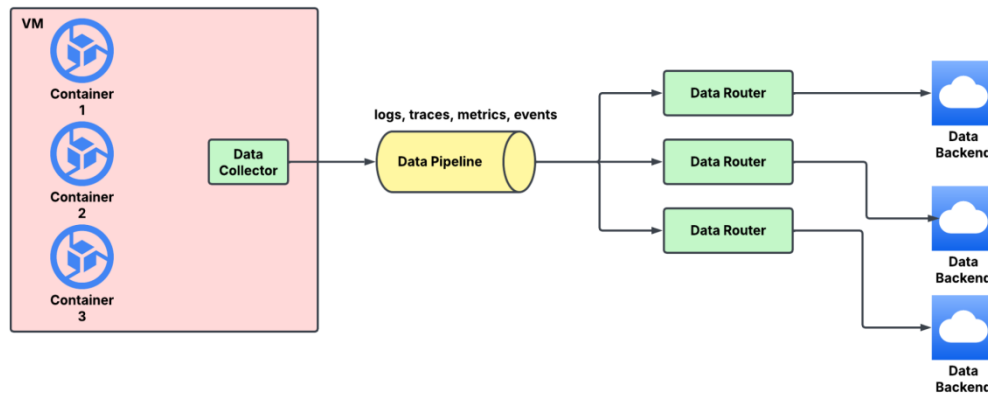


**FIG 2: A Layered Observability Framework Depicting Data Collection, Storage, And Analysis Pipelines**

**Proactive Incident Prevention and Faster Resolution**

While monitorability allows a team to predict incidents and address them before they turn into major events. Some of the use cases include basic and advanced metric visualization, key variances, and real-time anomaly detection for response times and error percentages [5]. Seven design principles include alerting mechanisms that make it possible for teams to receive immediate communication concerning possible problems, founded on the observability data.

Of all the use cases of observability, it quickens the rate of getting to the cause when solving incidents. For example, multiple metrics may show that a particular service is in poor condition while logs contain information about particular errors and traces to explain which precise service or endpoint triggered them. They achieved great satisfaction with this integrated view since it decreases the mean time to the resolution (MTTR) and losses.

**Observability Tools and Techniques**

Modern observability relies on a suite of tools and techniques designed for cloud-native environments:

*A. Prometheus:*

A post-discussion/service/outcome metering and tracking mechanism that uses a database to collect time-stamped figures and issues regular or on-demand alerts.

*B. Grafana:*

A dashboard tool that generates custom visualizations to interpret figures and patterns.

*C. Jaeger:*

A tool used to show request flows across numerous microservices and allow for latency and dependency analysis.

Together with the techniques for using the service anomalies, such as anomaly detection, predictive analytics, and automated alerting, these tools allow teams to handle the challenges of distributed systems proficiently [5].

## Comparison with Traditional Monitoring

Standard supervising relies on ratios and easily set TRIGGER levels giving poor insight into the systems. Observability enhances the collection of data throughout metrics, logs, and traces and then correlates this data, resulting in more insights. Another key difference is that monitoring is a passive process, whereas observability focuses on making bugs and other issues visible to the right people at the right time to increase system reliability and performance up-time.

## IV. DESIGNING AN OBSERVABILITY-DRIVEN INCIDENT MANAGEMENT FRAMEWORK

## Conceptual Framework for Observability-Driven Incident Management

Also, an observability-driven incident management framework is developed as a best practice approach for improving the reliability of cloud-native applications by incorporating observability techniques into each stage of the incident management process. The framework encompasses four key components: identification of events, data gathering, correlation as well as themselves, and automated mitigation. Together with these components arises the possibility of early issue identification; fast root-cause analysis; and the quick mitigation of incidents to prevent their adverse effects on the availability and performance of the system.

## Data Collection: Automated Logging and Telemetry

At the very base of the framework is the integration of data gathering and monitoring of distributed system components. As mentioned above, this environment means that metrics, logs, and traces in the form of telemetry data are collected continuously from containers, services and infrastructures [6]. Automated logging systems show the error messages and warnings and operation performed and telemetry gives information regarding the usage of resources such as CPU and network latency among others and transactions.

## Incident Detection: Real-Time Anomaly Detection Using Machine Learning

In the context of incident management within the proposed framework, the incident detection capability can be described as real-time anomaly detection using machine learning models. These models process telemetry information to look for unusual behaviors from a subject's baseline profile. For example, excessive response time, high error rate or CPU usage is considered to be an incident that has occurred. Clustering and actuarial modeling further improve the effectiveness of the recognition of deviations by taking into consideration temporal cycles, fluctuations in system loads, and service-specific benchmarks. Inter-connectivity with real-time alerting systems means that deviations result in alerts for the right people so that action can be taken quickly.
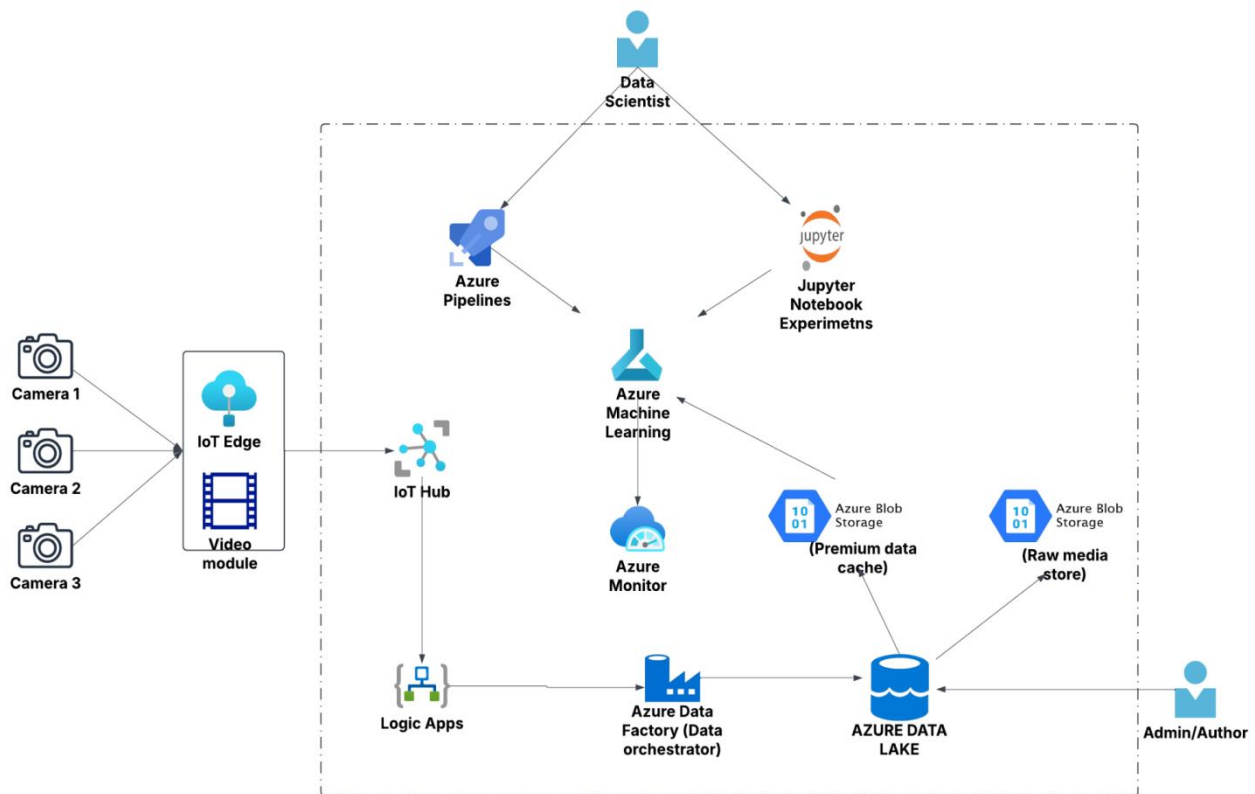
**FIG 3: A System Architecture Design for Observability-Driven Incident Management, Including Data Flow and Decision-Making Components**

## Correlation Analysis: Tracing Incidents to Specific Services or Nodes

After an incident occurs, correlation analysis is done by the framework to trace the problem to its source. In this phase, the distributed tracing tools, like Jaeger or OpenTelemetry, as explained above, help in mapping the flow of the requests and identifying the service or the node that is actually serving the request.

Correlation analysis involves the use of a combination of information components to come up with a coherent picture of an event. For example, if analysis of throughput indicates that latency is high, error messages found in logs may provide association of such problems to certain traces, and these may show the full sequence of transactions made by the transaction [7].

## Automated Remediation: Deployment of Self-Healing Mechanisms

The framework employs the automated remedial measures to prevent system unavailability and to ensure system stability. The actions like service restart, adding or reducing resources and rerouting the traffic is from predefined policies. For instance, if a service is observed to contain high CPU usage, then, the framework will create more copies to accommodate the usage.

Derivative orchestration software, including Kubernetes, offers support for implementing automated remediation by managing the services expressed as containers. Interaction with other observability tools guarantees that remediation procedures are intelligent and relevant to the environment.

## Challenges in Implementing the Framework

Implementing an observability-driven incident management framework presents several challenges:

A. *Data Storage and Processing Overhead:*

The high volume of telemetry data gathered and its processing operational loads on devices involves large investments. This rapid growth presents unique issues such as data compression, data retention policies, and the general scalability of cloud storage solutions to support this data doubling every two to five years in enterprises [8].

*B. False Positives:*

It is possible, in some cases, that the first type of machine learning models used for anomaly detection might give out a considerable number of false positives which in turn result in alerts and rectification actions. In order to avoid this kind of problem, constant model calibration and threshold tuning are crucial.

*C. Integration Complexity:*

While it is possible to bring observability tools into the incident management system, mostly it is quite challenging and engineering heavy work, especially in complex and diverse setups.

*D. Performance Impact:*

Real-time acquisition and analyses can add new latency to processing that can be detrimental to performance, especially in limited resource settings. The issue of how observable a system should be, as well as how well it should perform, is a critical factor to address.

**Potential Benefits of the Framework**

Despite its challenges, an observability-driven incident management framework offers significant benefits:

*A. Reduced Mean Time to Resolution (MTTR):*

The main aim of the framework presented here is to deliver tangible recommendations and automate remediation processes, which can save considerable time on recognizing threats and investigating, as well as remediating incidents.

*B. Improved System Uptime:*

Incorporation of the PA Hunting with self-correction mechanisms will enable better system availability, minimal downtime and overall system reliability [9].

*C. Enhanced Developer Productivity:*

Managers and developers get rid of time-wasting incidents by having well-organized mechanisms to track and handle incidents and concentrate on delivering new features and advancing their product.

*D. Scalability and Adaptability:*

The framework specifically aligns with the fast-paced, dispersed architecture that is common in cloud-native systems to guarantee stable results at scale.

## V.  CASE STUDIES AND APPLICATIONS

**Example 1: Scaling Microservices Architecture with Real-Time Performance Monitoring**

A real-life example can be derived from a well-established e-commerce platform that experienced issues with the scale of its microservices on Black Friday. Performance management was done with an observability-first approach with the aid of tools which includes Prometheus as with Grafana. Actual real-

time monitors were call rates per second, CPU load, and memory consumption. It allowed the team to reduce high CPU latency in payment processing services during congestive time. Jaeger was also used to identify certain microservices that were experiencing high latency making optimization easier [10]. Equipped with this data, they used several scaling techniques, including the deployment of certain heavy computational services across several hosts, and also enhancement of caching practices to deal with as few database queries as possible. They achieved these improvements through optimization leading to 40% improvement in response time and guarantee of quality during congestion.

## Example 2: Automating Incident Resolution in a Kubernetes Cluster

A company offering financial services deployed observability for solving incidents in Kubernetes infrastructure. The collected logs, metrics, and traces using the chosen instrument – Open Telemetry, which structures the information from the containerized applications and consolidates it in the control center. The information was then analysed using machine learning algorithms which was actively checking for irregularities. For instance, large variances in error rates warned other groups, with traces linked to said alerts to pinpoint affected services. According to predefined remediation policies the corrective actions such as restarting pods, scaling deployments or even traffic redirecting were performed [9].

Memory leak was evidenced by one of the critical services where pod crashed and in effect, disrupted the operations. Thanks to the observability framework, the anomaly was recognized, isolated at the level of a particular service, and an automatic corrective process was initiated to restart the pod and scale the additional number of replicas. Thanks to this fast decision, the company lost little time, and customers also suffered little or no consequences. This has suggested that the implementation of the monitoring and response process has achieved the following improvements: average time to response and resolution cut 50% and the use of manual interventions that cost 30% less than expected.

## Results: Improved Reliability, Cost Savings, and User Satisfaction

Reflecting the practice, the case studies focus on proving the effectiveness of analysis-driven strategies in terms of system availability, service cost, and users' satisfaction. Various stakeholders of information systems realized improved system availability and minimized disruptions by leveraging real-time insights into system activities and sponsoring auto-generated response for managing incidents. For example, the platform's assignment slowness-related with the overloaded traffic which is negatively affecting the customer's satisfaction and sales conversion. There were also savings; automated remediation was made possible hence decreasing on operations teams having to manually fix issues.

## Limitations and Lessons Learned

Despite the numerous benefits offered by observability-driven approaches, several challenges have been identified. One concern was that telemetry data coming from microservices architectures were becoming almost overwhelming to be processed [12]. Companies soon began to understand the need for implementing better ways of data collection and storage as a response to these issues. Other problems were cases of missing valid information in certain anomalies or false positive cases that invoked numerous alerts and autos. Furthermore, it was emphasized that adjusting the detection models with reference to regular intervals and setting up reference points were indispensable steps of the action plan to stabilize the system.

## VI. CONCLUSION

This report underscores the three pillars of this architecture, with observability being key to credibility of cloud-native applications. In today's distributed environments nomenclature techniques are inadequate

when coping with situations of slowness, steering, and overall rewrite. When presented as metrics, logs, and traces, observability allows you to manage the systems' behavior proactively, detect incidents, analyze the root causes accurately, and resolve issues quickly. Increased observability is a key component of the presented structure for implementing an effective incident management approach that includes real-time data gathering, anomaly identification, correlation, and automatic resolving to improve performance. When these practices are adopted, it becomes possible for organizations to reduce mean time to resolution (MTTR), increase the system uptime and increase on the productivity of the developers. This framework not only accelerates issue resolving but also allows teams to grow and seamlessly deploy the cloud-native systems. As the ways in which AI generates information become more observable in the future, handling incidents will again be affected. There should be increased alerts from machines on the anomalies which should do away with lots of alarms with improved specificity; There should also be better chances of automated rectification by the system on the anomalies.

## REFERENCES

[1] D. Gannon, R. Barga, and N. Sundaresan, "Cloud-Native applications," *IEEE Cloud Computing*, vol. 4, no. 5, pp. 16–21, Sep. 2017.

[2] Chakraborty, Mainak & Kundan, Ajit. (2021). Monitoring Cloud-Native Applications: Lead Agile Operations Confidently Using Open Source Software. 10.1007/978-1-4842-6888-9.

[3] R. Picoreti, A. P. D. Carmo, F. M. De Queiroz, A. S. Garcia, R. F. Vassallo, and D. Simeonidou, "Multilevel Observability in Cloud Orchestration," *Multilevel Observability in Cloud Orchestration*, pp. 776–784, Aug. 2018.

[4] N. Marie-Magdelaine and T. Ahmed, "Proactive Autoscaling for Cloud-Native Applications using Machine Learning," GLOBECOM 2020 - 2020 IEEE Global Communications Conference, Taipei, Taiwan, 2020, pp. 1-7, doi: 10.1109/GLOBECOM42002.2020.9322147.

[5] M. Wurster, U. Breitenbücher, A. Brogi, F. Leymann, and J. Soldani, "Cloud-native Deploy-ability: An Analysis of Required Features of Deployment Technologies to Deploy Arbitrary Cloud-native Applications," *Cloud-native Deploy-ability: An Analysis of Required Features of Deployment Technologies to Deploy Arbitrary Cloud-native Applications.*, pp. 171–180, Jan. 2020.

[6] R. Peinl, F. Holzschuher, and F. Pfitzer, "Docker Cluster Management for the cloud - survey results and own solution," *Journal of Grid Computing*, vol. 14, no. 2, pp. 265–282, Apr. 2016.

[7] J. Kosińska and K. Zieliński, "Autonomic Management Framework for Cloud-Native Applications," *Journal of Grid Computing*, vol. 18, no. 4, pp. 779–796, Sep. 2020.

[8] Q. Duan, "Intelligent and Autonomous Management in Cloud-Native Future Networks—A Survey on Related Standards from an Architectural Perspective," *Future Internet*, vol. 13, no. 2, p. 42, Feb. 2021, https://www.mdpi.com/1999-5903/13/2/42.

[9] K. A. Torkura, M. I. H. Sukmana, F. Cheng, and C. Meinel, "Leveraging Cloud Native Design Patterns for Security-as-a-Service Applications," *Leveraging Cloud Native Design Patterns for Security-as-a-service Applications*, pp. 90–97, Nov. 2017.

[10]   A. Al-Fuqaha, M. Guizani, M. Mohammadi, M. Aledhari, and M. Ayyash, "Internet of Things: A Survey on Enabling Technologies, Protocols, and Applications," *IEEE Communications Surveys & Tutorials*, vol. 17, no. 4, pp. 2347–2376, Jan. 2015.

[11]   M. Dickinson *et al.*, "Multi-Cloud performance and security driven federated workflow management," *IEEE Transactions on Cloud Computing*, vol. 9, no. 1, pp. 240–257, Jun. 2018.

[12]   K. A. Torkura, M. I. H. Sukmana, F. Cheng, and C. Meinel, "Leveraging Cloud Native Design Patterns for Security-as-a-Service Applications," *Leveraging Cloud Native Design Patterns for Security-as-a-service Applications*, pp. 90–97, Nov. 2017.

[13]    N. Kratzke and R. Siegfried, "Towards cloud-native simulations – lessons learned from the front-line of cloud computing," *The Journal of Defense Modeling and Simulation Applications Methodology Technology*, vol. 18, no. 1, pp. 39–58, Jan. 2020.