

Secure Kubernetes Reference Architecture for Enterprise and Artificial Intelligence Workloads

Navaneeth Komirisetty

Sr.Cybersecurity Architect at American Express Travel Related Services inc., USA

Abstract

The shift to enterprise Kubernetes has transformed the responsibility of platform teams from cluster operators to overseers of multi-tenancy products with security-by-default policies. AI workloads alongside application services bring high-value assets with unique security properties to the platform, such as training datasets, model artifacts, and inference endpoints. We present a vendor-neutral reference architecture for Kubernetes-based platforms for traditional workloads and AI. Using design science research methodology, we integrate known security frameworks, research on compliance automation, and findings from practitioner literature. Trust boundary segmentation, tenant isolation, control plane hardening, policy-based workload control, and evidence-ready observability form the core components. In particular, it accounts for inference security for artificial intelligence pipelines and model-serving services, including verification of artifact provenance and segmentation of access to training and evaluation datasets. Architecture components are verified against industry security benchmarks and recognized governance frameworks. A progressive roadmap is defined, from baseline security controls to operation in a mature, continuously validated state across multiple teams and environments. The contribution extends Kubernetes security documentation, addressing emerging needs for enterprise governance of artificial intelligence workloads deployed on Kubernetes.

Keywords: Kubernetes Security, Enterprise Container Platform, Multi-Tenant Architecture, Cloud-Native Modernization, Zero Trust, Artificial Intelligence Workloads, Supply Chain Security

1. Introduction

1.1 Background and Context

Enterprise container platforms have become the underlying execution layer for modern digital services. Kubernetes is the cluster orchestration system of record for enterprise workloads in containers. Kubernetes has evolved over a decade from earlier enterprise-grade cluster management systems that focused on deploying, scaling, and managing workloads. Research on container management systems Borg, Omega, and Kubernetes identifies several tradeoffs, such as efficiency, fault tolerance, ease of operation, and security isolation (Burns et al., 2016). Containerization provides consistent runtime environments and simplified dependency management, which introduces new security considerations requiring systematic and architectural approaches to security.

Container implementations are empirically shown to offer a substantial improvement in start-up time, measuring at 50 milliseconds compared to virtual machine initialization taking 30-40 seconds (Sultan et al., 2019), as well as reduced kernel overhead by virtue of the shared kernel, which has led to significant demand for container products in the market. In 2016, the market for container products was 762 million US dollars and it is expected to increase to 2.7 billion US dollars by 2020 (Sultan et al., 2019), which makes it essential to develop security architectures to keep pace with the rapid adoption of container technology.

Kubernetes is also a frequently studied phenomenon in academic papers and reports with respect to deployment patterns, security issues, and operational challenges. Shamim et al. (2022) performed a literature review of 321 artifacts and 105 peer-reviewed research papers on Kubernetes and found it to be a frequently researched topic with respect to security, one of the most commonly discussed challenge categories in Kubernetes deployment.

1.2 Problem Statement and Research Gap

Existing security documents for Kubernetes-based platforms offer general hardening principles and apply standard practices to AI workloads, treating them in effect as any other service. There are, however, prominent differences between AI workloads and standard services. Training data may be subject to jurisdictional restrictions on PII. Model weights tend to be costly to compute and have business value, making them susceptible to exfiltration. Inference

endpoints support untrusted inputs at scale and can be exploited for prompt injections, data extraction, and resource consumption.

The gap is also quantitatively evidenced in practitioner literature. An analysis of 321 Kubernetes-related artifacts found the lack of security practices and tools to be the most commonly reported challenge, appearing in 41.1% of reviewed artifacts (Shamim et al., 2022). The top two concerns by practitioners were attack surface reduction (39.2%) and lack of diagnostics tools (37.7%). Only 6.7% of the 105 peer-reviewed publications focused on security aspects of APIs. This may indicate a gap between the interests of practitioners and research literature (Shamim et al., 2022). In addition, it is notable that there are no research publications for nine of the 15 challenges.

Container security research has identified classes of vulnerabilities that require architectural mitigations, and container-based deployments in cloud computing environments have also been shown to experience isolation issues in shared kernel architectures compared to conventional virtualization approaches (Tak et al., 2017). When conducting multiple experiments to identify vulnerabilities in container security, the default configurations of 56.8 percent of container images were vulnerable (Sultan et al., 2019). Based on the National Institute of Standards and Technology Artificial Intelligence Risk Management Framework, AI systems can be subjected to existing vulnerabilities in traditional cybersecurity, as well as new vulnerabilities from machine learning models (NIST, 2023). These properties require security architectures beyond standard application protection patterns.

1.3 Contribution and Scope

This article presents a reference architecture for both traditional workloads and artificial intelligence on enterprise Kubernetes. It describes an architecture that combines practices for container security with new practices for artificial intelligence governance. The architecture emphasizes principles that enterprises can implement with their existing capabilities, rather than specifying particular commercial products.

It addresses the research gap reported in previous work in three ways: first, through a comprehensive literature review of academic and practitioner sources for operational relevance using the multi-vocal approach; second, by leveraging empirical findings from compliance automation research to demonstrate security benefits for the proposed framework; third, this framework accounts for a set of artificial intelligence workload profiles not considered in existing Kubernetes security recommendations.

This includes multi-tenant platform design principles, control plane security fundamentals, node security fundamentals, workload policy, observability for security operations, and patterns for protecting artificial intelligence workloads. An implementation roadmap provides guidance for how the architecture principles translate into deployment recommendations. The intended audience includes security architects, platform engineers, cloud infrastructure leads, and security operations professionals responsible for the secure design and operation of multi-tenant container platforms.

2. Related Work

2.1 Container Security Frameworks

As containerization has grown, so has research into container security. Container security is broken down into five broad risk areas outlined by the National Institute of Standards and Technology Application Container Security Guide: image, registry, orchestrator, container, and host operating systems risks (Souppaya et al., 2017). This classification system has been widely adopted in both research literature and practitioner guidance, providing a common language to discuss container security issues.

An analysis of container security issues classified the host-container threat model into four use cases: protecting a container from applications inside it, inter-container protection, protecting a host from containers, and protecting containers from malicious or semi-honest hosts (Sultan et al., 2019). The first three use cases are mostly addressed by software solutions that rely on Linux kernel features such as namespaces, cgroups, capabilities, and seccomp. The fourth use case is addressed by hardware solutions such as trusted platform modules (TPMs) and Intel SGX. This classification is useful to describe security controls but does not capture the specific requirements of AI workloads.

2.2 Kubernetes Security Research

Although Kubernetes has widespread practical usage, only a limited number of research papers have studied Kubernetes security from the practitioners' perspective. One literature review identified seven out of 105 peer-reviewed papers as relating to security concerns (Shamim et al., 2022). The identified security concerns include anomaly detection, automated response to security incidents, and zero trust secure design in Kubernetes environments. However, none of the documents reviewed address the security of artificial intelligence workloads.

Despite the availability of these best practices, surveys published in the industry have shown that practitioners are concerned about the security of Kubernetes. In one survey, 44 percent of respondents indicated that they delayed deploying Kubernetes-based software due to lack of confidence in Kubernetes security (Shamim et al., 2022). Another survey indicated that 67 percent of respondents misconfigured their Kubernetes deployments in a manner that manifested security issues (Shamim et al., 2022). This suggests security concerns are significant barriers to widespread adoption and require systematic architectural solutions.

2.3 Compliance Automation Research

A study on automating CIS Benchmark compliance in Kubernetes worker nodes found that automation can be 99.8 percent faster than manual compliance, with authors noting that automated compliance can reduce the cost to an organization (Sharma et al., 2024). The same study showed that automated systems reduced persistent benchmark failures by 95 percent and resulted in a 99.5 percent compliance posture during continuously running scans, validating the policy-driven enforcement approaches adopted in the reference architecture.

2.4 Artificial Intelligence Governance

Several frameworks have been proposed for governing widely deployed artificial intelligence systems, including the NIST Artificial Intelligence Risk Management Framework, which provides a structured approach for identification and categorization of artificial intelligence risks across the system lifecycle (NIST, 2023). The Generative Artificial Intelligence Profile builds on this framework to identify risks associated with large language models and generative systems (NIST, 2024). However, while these frameworks address artificial intelligence governance at the organizational level, they provide limited guidance at the platform level in containerized environments.

2.5 Differentiation of this Work

In contrast to the existing body of literature, this article integrates existing research on container security with the requirements of artificial intelligence governance, which have mainly developed as separate streams. Second, it provides empirical contributions based on existing scholarship and practitioner literature through multi-vocal analysis to ensure operational relevance. Third, it provides implementation-oriented contributions through a phased roadmap and validation framework to complement conceptual architecture with implementation and deployment patterns.

3. Research Methodology

3.1 Design Science Research Methodology

The research followed design science research methodology, which is tailored to IS artifacts such as reference architectures (Hevner et al., 2004). Design science research is concerned with the creation and evaluation of artifacts intended to solve identified organizational problems. The approach consists of problem definition, goal specification, design and development, demonstration, evaluation, and communication.

The problem identification phase discussed in Section 1.2 focused on the analysis of previous work and practitioner reports to identify the research gap. The goal definition phase focused on establishing design requirements concerning the support of traditional and artificial intelligence workloads, vendor neutrality, established security frameworks, and implementability by practitioners. The reference architecture was built iteratively during design and development, instantiated from existing architectural frameworks, empirical results, and architectural patterns.

3.2 Knowledge Base Integration

Knowledge from different subject areas is integrated based on principles of design science. The architecture relies on NIST and CIS cybersecurity frameworks and on zero trust architecture. Empirical knowledge is based on multi-vocal literature review findings that reflect practitioners' challenges and research gaps (Shamim et al., 2022). Technical

knowledge is based on documented container security practices provided by systematic analysis (Sultan et al., 2019; Souppaya et al., 2017). Compliance knowledge is based on automation research, which has provided measures for improving security (Sharma et al., 2024).

3.3 Design Principles

The architecture is based on five guiding principles derived from the knowledge base.

To address findings that 56.8 percent of exploits could succeed against default containers, defense in depth proposes that security controls should be implemented at various layers of system architecture, rather than at a single security boundary (Sultan et al., 2019).

Policy-driven enforcement mandates automation to achieve policy compliance. Empirical studies have shown that policy-driven enforcement reduces the time for compliance by 99.8 percent when compared to manual enforcement (Sharma et al., 2024).

Zero trust boundary definition does not consider trust to be based on physical or logical location, or on identity claims inside the boundary. This aligns with research defining comprehensive visibility over assets and continuous monitoring as requirements of successful zero trust adoption (He et al., 2022).

Evidence-ready operations should create audit evidence during normal operations, rather than through separate compliance activities, addressing practitioner concerns about compliance burden noted in multi-vocal literature (Shamim et al., 2022).

Artificial intelligence asset protection: assets such as training data, model artifacts, and inference endpoints must be clearly identified and protected as first-class assets. This principle addresses the gap between general container security recommendations and artificial intelligence-specific recommendations of governance frameworks (NIST, 2023; NIST, 2024).

3.4 Evaluation Approach

An expert review assesses whether the architecture meets the criteria, which include security coverage, compliance with established frameworks, and implementability (Section 11). The evaluation framework addresses peer review concerns regarding the validation of proposed architectural components.

4. Threat Model And Trust Boundary Analysis

4.1 Attack Surface Characterization

Secure platform design starts with threat modeling of the assets being protected and adversary capabilities. Kubernetes attack surfaces span multiple layers in the platform architecture, which should be addressed individually. The National Institute of Standards and Technology Application Container Security Guide identifies five high-level risks in containerized deployments: image risks, registry risks, orchestrator risks, container risks, and host operating system risks (Souppaya et al., 2017). Each of these risks can manifest in multiple types of vulnerabilities.

Image risks include base image vulnerabilities, build-time code injections, and improperly handled secrets (e.g., saving credentials to layers). Registry risks include poor access control, lack of image verification, and absence of vulnerability scanning. Risks to the orchestrator include weak authentication and authorization, and lack of audit logs. Risks to the containers include privilege escalation, container escape, and failure of network segmentation. Host risks include vulnerabilities in the kernel, patching issues, and granting excessive privileges for containers to access resources on the host (Souppaya et al., 2017).

In one study of container security concerns, the attack surface was found to be larger due to containers sharing the host kernel (Sultan et al., 2019). In contrast to containers, traditional hypervisor-based isolation boundaries are stronger but suffer from tradeoffs with respect to resource efficiency and security assurance. The evaluation study further noted that 56.8 percent of the evaluated exploits succeeded against default container images, highlighting the need for defense-in-depth protection based on multiple methods of isolation, rather than reliance on a single security boundary (Sultan et al., 2019).

4.2 Artificial Intelligence Asset Classification

Asset types include artificial intelligence workloads that are distinct from application data. The NIST Artificial Intelligence Risk Management Framework provides a taxonomy to identify and categorize artificial intelligence-related risks throughout the system lifecycle (NIST, 2023). The framework distinguishes between risks to people, organizations, and ecosystems, in which artificial intelligence systems can cause different levels of harm than conventional software systems.

Training datasets may include personal or proprietary business information, as well as licensed third-party intellectual property that must be used in accordance with contractual stipulations. The Generative Artificial Intelligence Profile expands the Baseline Risk Management Framework (RMF) to include risks specific to large language models and other generative AI systems (NIST, 2024). These risks include privacy risks from leaking information from the training set, intellectual property risks from accidentally generating copyrighted material, and integrity risks from generating false or misleading information.

Challenges from the Generative Artificial Intelligence Profile also include model integrity. Model weights, which represent the majority of the computational output of the training process, represent meaningful investments for the organization; these could be altered by an adversary to produce subtly wrong and harmful outputs in a way that can evade detection (NIST, 2024). This covers risk scenarios beyond traditional data integrity challenges that can be attributed to behaviors of machine learning systems that static analysis cannot detect.

4.3 Trust Boundary Definition

A trust boundary is a boundary of the environment with different security and verification requirements. Zero trust architecture principles guide the definition of such boundaries. A study investigating zero trust implementation barriers found that a complete transition to zero trust requires comprehensive visibility of network traffic, strong identity verification mechanisms, and continuous monitoring and detection of access behavior (He et al., 2022). Another study found that legacy system integration, performance overhead, and lack of organizational buy-in for fundamental access control philosophy change are the main barriers to zero trust adoption.

One reference architecture employed by many organizations for establishing a zero trust architecture has five pillars: identity governance, device security, network segmentation, application security, and data classification (Phiayura & Teerakanok, 2023). In terms of network security, micro-segmentation alone is not sufficient without enhancement to identity validation and continuous authorization. The reference architecture contains explicit boundaries: external boundaries separate public networks from platform ingress; administrative boundaries separate platform operators from tenant workloads; tenant boundaries prevent lateral movement between tenants; workload boundaries separate deployments within a tenant namespace from one another.

Asset Category	Examples	Protection Requirements	Threat Vectors
Training Data	Datasets, feature stores	Encryption, access control, lineage tracking	Theft, poisoning, compliance violations
Model Artifacts	Weights, embeddings, indexes	Integrity verification, provenance, versioning	Exfiltration, tampering, unauthorized access
Inference Endpoints	Model APIs, serving infrastructure	Input validation, rate limiting, monitoring	Prompt injection, data extraction, denial of service
Configuration	Prompts, templates, parameters	Version control, access restriction, audit	Manipulation, unauthorized modification
Operational Data	Logs, metrics, traces	Retention, access control, tamper protection	Evidence destruction, privacy exposure

Table 1: Asset Classification and Protection Requirements [NIST (2023); NIST (2024); Souppaya et al. (2017)]

5. Multi-Tenant Isolation Architecture

5.1 Tenancy Model Design

Multi-tenancy also determines the blast radius when a compromise occurs, making the isolation model one of the most critical architectural decisions in the cloud. Isolation models must consider the tradeoffs between resource usage and security. More isolation means more infrastructural and operational overhead. Container platforms provide a variety of isolation solutions that can be combined to achieve security.

An empirical study of Kubernetes adoption has found multi-tenancy to be one of the main challenges for enterprise organizations (Shamim et al., 2022). Multi-vocal literature review found Kubernetes security literature heavily focused on resource sharing, access control, and workload isolation. Learning curve challenges, affecting 28.3 percent of 321 internet artifacts, along with complex Kubernetes concepts, impede isolation configuration (Shamim et al., 2022).

Container isolation is achieved using Linux kernel features like namespaces, cgroups, and capabilities. A 2017 paper on the security of containers in the cloud stated that namespaces provide logical isolation but do not eliminate all forms of cross-tenant interference (Tak et al., 2017). A cluster's kernel components such as its network stack, file system, and system call interface may also lead to increased information leakage (Sultan et al., 2019). Research on the design of cluster management systems found that multiple layers of isolation for workloads in virtualized environments are necessary to meet reliability and security requirements (Burns et al., 2016).

5.2 Workload Classification Framework

For effective isolation, workload classification according to data sensitivity, regulatory requirements, and attack surface should be established. The NIST Artificial Intelligence Risk Management Framework recommends risk-based approaches to select security controls with intensity proportional to potential impact of system compromise or failure (NIST, 2023). This can help identify isolation tiers, monitoring levels, and access control standards.

Container security guidance typically begins by defining the containers' workload privileges (Souppaya et al., 2017). The Application Container Security Guide defines different workloads based on whether they require host access, high container privileges, or whether they can run in a restricted execution context. Classification directly influences isolation tier assignment and policy application to workloads.

Taxonomies should include issues related to traditional software application risks as well as risks specific to AI systems. The Generative Artificial Intelligence Profile outlines categories of risks associated with the deployment of large language models, including the sensitivity of the training data, the model's capabilities, and integration with third-party systems (NIST, 2024). A production inference endpoint that processes personal data may be regulated-tier, whereas experimental training jobs on synthetic data may be in standard isolation tiers.

5.3 Identity and Access Governance

To implement isolation management effectively, strong identity management is needed for humans and machines. Zero trust architecture principles help improve isolation by continually verifying claims about human and machine identities, instead of trusting their network location as a means for authorization (Phiayura & Teerakanok, 2023). Identity governance is a foundation that needs to be in place before other zero trust efforts.

A survey on zero trust deployment highlighted the need to accommodate both interactive user sessions and programmatic service-to-service interactions in identity validation (He et al., 2022). Service identity is especially difficult in containerized environments where workloads are dynamic and ephemeral. Workload identity solutions require cryptographic identity attestation and operational flexibility.

Network isolation should complement identity-based access controls by restricting communication regardless of identity claims. Network policies should default-deny traffic unless explicitly allow-listed, matching specific known communication paths for further isolation. Artificial intelligence inference endpoints should have egress policies limiting model-serving containers from opening unexpected outbound connections.

Tier	Isolation Level	Appropriate Workloads	Control Requirements
Standard	Namespace	Internal tools, development, non-sensitive batch	Baseline pod security, network policies, resource quotas
Elevated	Namespace with dedicated node pools	Customer-facing services, confidential business data	Enhanced monitoring, stricter policies, access reviews
Regulated	Dedicated node pools or clusters	Compliance-mandated processing, jurisdictional data	Documented controls, audit evidence, external attestation
AI-Specific	Dedicated accelerator pools with namespaces	Model training, inference serving, data pipelines	Data access segmentation, artifact governance, specialized monitoring

Table 2: Isolation Tier Characteristics [Souppaya et al. (2017); Shamim et al. (2022); Tak et al. (2017); Phiayura & Teerakanok (2023)]

6. Control Plane and Node Security Foundation

6.1 Control Plane Hardening

The Kubernetes control plane is the security control mechanism of the Kubernetes platform. If the control plane is compromised, an attacker can take over the entire cluster, alter workloads, retrieve secrets, and implant backdoors. Thus, orchestrator security is a separate category of threats that requires specific mitigations (Souppaya et al., 2017). Control plane hardening involves access control, auditability, data protection, and operational controls.

The largest attack surface within the control plane is access to the API server. Research into container orchestration security has shown that an attacker with access to the API server has control over the entire cluster (Sultan et al., 2019). Common patterns of production deployments potentially exploited by attackers include authentication issues, authorization issues, and insufficient audit logging. The recommended best practices are multilayered defense, strong authentication, fine-grained authorization, and comprehensive security logging.

Audit logging should capture all API activity to a sufficient level of detail to detect and investigate incidents and support compliance activities. Container security recommendations advise that logs be stored in security information systems instead of being stored on control plane nodes where they are accessible to an adversary for deletion (Souppaya et al., 2017). Sensitive audit events (e.g., secret access, privilege escalation, administrative commands) should be logged and alerted to the security operations team.

6.2 Worker Node Hardening

Worker nodes that run the actual tenant workloads should be hardened to provide the same level of security as the tenant applications. The Center for Internet Security Kubernetes Benchmark describes hardening for node configuration, kubelet, and container runtime. Given the size of enterprise clusters, which can contain hundreds or thousands of worker nodes, manual compliance checking is not feasible (Sharma et al., 2024).

An automated compliance assessment of Kubernetes worker nodes found that achieving CIS Benchmark compliance is best served by continuous monitoring rather than point-in-time assessment, lowering mean time-to-compliance from over 120 minutes with manual remediation to under two minutes when using automated remediation (Sharma et al., 2024). The compliance posture was 99.5 percent during continuous load testing, but it would drop to 78 percent without automation. Configuration drift, in which nodes would no longer be compliant after software updates and/or administrative or workload changes, was also a problem.

Multiple complementary controls should be used to reduce the container runtime attack surface. The Application Container Security Guide recommends seccomp profiles to limit system calls in order to reduce impact of container escape (Souppaya et al., 2017). Research on system call filtering found that automated profiling reduces the number of available system calls by more than 35 percent without any meaningful performance cost (Sultan et al., 2019), while security namespace implementation incurs less than 0.7 percent overhead in latency (Sultan et al., 2019).

6.3 Accelerator Node Considerations

AI workloads often use hardware accelerators to improve performance, including graphics processing units and tensor processing units. This adds additional security concerns to the workloads, as the accelerator driver and runtime libraries now become part of the trusted computing base (TCB), which must be maintained as controlled artifacts with documented provenance and timely updates. Containerization in the cloud has shown that device access from inside containers raises security issues that do not exist in CPU-only workloads (Tak et al., 2017).

Device plug-ins exposing accelerators to Kubernetes pods should ensure that they are configured to only be available to authorized workloads. Dedicated node pools for nodes with accelerators are an effective way to reduce the amount of configuration required, since node affinity rules can be used to select appropriate infrastructure for artificial intelligence workloads. A taint and toleration mechanism prevents non-artificial intelligence workloads from scheduling on accelerator nodes, preserving the availability of these resources and reducing the attack surface.

According to the NIST Artificial Intelligence Risk Management Framework, the security of an AI system includes the infrastructure on which model training and inference are performed (NIST, 2023). Security features at the hardware level include secure enclaves or memory encryption where high-sensitivity workloads handle sensitive data or proprietary model weights.

7. Policy-Driven Workload Governance

7.1 Pod Security Standards Implementation

Workload policy enforcement is most effective when done automatically, rather than manually or by convention. The Application Container Security Guide recommends a baseline set of policies for all application workloads, limiting privilege escalation, Linux capabilities, and access to host resources (Souppaya et al., 2017). The Kubernetes Pod Security Standards provide a taxonomy that enforces workload privilege levels at the namespace or cluster level.

Research into container security issues identified privilege escalation as a major risk to reduce (Sultan et al., 2019), as containers with higher privileges can escape their isolation and access the host and other containers. Policy enforcement at the admission layer prevents non-compliant workloads from being deployed before security violations are created.

Implementation is phased with an initial audit mode, wherein violations are logged without blocking workload deployment. Exceptions to baseline policies should be documented, limited in duration, and accompanied by additional monitoring requirements, as stated in the Application Container Security Guide (Souppaya et al., 2017). Organizations should create timelines that meet security improvement needs with minimal operational disruption.

7.2 Admission Control and Policy Enforcement

Pod Security Standards define pod-level security properties, while additional governance requirements define custom admission control policies. Admission controllers are software components that intercept requests to create or modify workloads in a cluster and evaluate them against organizational policies before allowing changes to the cluster's state. Admission-time enforcement offers earlier visibility than runtime detection in the context of continuous security assessment of microservices environments (Torkura et al., 2017).

Container image provenance verification can ensure that container images are only executed if they have been signed by trusted registries. The Application Container Security Guide lists image integrity as a foundational security requirement, recommending cryptographic verification of image contents before deployment (Souppaya et al., 2017). Beyond signing images, supply chain security includes the selection and provenance of base images, dependencies, and the build process.

An analysis of production use of Kubernetes found that resource management was one of the main reliability and security problems (Vayghan et al., 2021). Without resource limits, workloads can over-consume resources in the cluster and effectively perform a denial of service against other tenants. Resource requests and limits also allow for capacity planning and cost allocation. Research into compliance automation found that 80% of the 65 CIS controls relevant to worker nodes could be applied with remediation and admission control (Sharma et al., 2024).

7.3 Secret and Configuration Management

Secrets management is an important control for workloads that use model registries, dataset storage, external application programming interfaces, and service credentials. Souppaya et al. (2017) emphasize in the Application Container Security

Guide that secrets should be deployed to separate secrets management systems, rather than Kubernetes secrets objects. Any workload access should be provided via service account bindings, not environment variable injection.

Zero trust architecture research established ephemeral workload identity lifecycle management as a core component (Phiayura & Teerakanok, 2023). Secrets rotation is a process with continuous improvement, not a one-time compliance exercise. Regular automated service account and managed identity credential rotation reduces the exposure window if a credential is leaked or stolen.

Environment-specific data should be read from an external configuration store or Kubernetes config maps. Kubernetes manifests should be managed using an infrastructure as code approach, allowing for code reviews, versioning and rollback of configuration changes, and for the history of changes to be auditable for compliance purposes.

8. Observability and Security Operations

8.1 Security Monitoring Architecture

Security architectures depend on operational visibility for timely detection and response. Security research on continuous security analysis of cloud-native applications has shown that effective security monitoring depends on fusing information from multiple sources of telemetry data (Torkura et al., 2017). Siloed monitoring tools provide only limited visibility to attacks, while integrated tools can detect multi-stage attacks across application and infrastructure layers.

Kubernetes audit logs track API requests and the associated authentication, authorization, and resource interactions. Kubernetes audit logging is covered in the Application Container Security Guide, which defines audit logging as a foundational security control enabling real-time detection and forensic investigation (Souppaya et al., 2017). The audit policy should be configured to capture a sufficient level of detail for security-relevant events without generating excessive logs for high-frequency operations.

Cybersecurity operations research states security monitoring has value through the presence of technical capabilities and organizational processes (Kovács, 2018). Security monitoring has value when paired with incident response capabilities to act on detected threats. Platform architectures should focus on the entire detect-and-respond cycle rather than merely telemetry collection.

8.2 Artificial Intelligence Workload Monitoring

The observability of AI services requires additional measurements beyond standard applications. The NIST Artificial Intelligence Risk Management Framework identifies monitoring as a key activity in risk management across the artificial intelligence system lifecycle (NIST, 2023). Model-serving endpoints need to expose request-volume metrics, latency distributions, error rates, and resource consumption metrics to enable anomaly detection.

A survey of methods for monitoring the behavior of artificial intelligence systems highlighted that the behavior of large language models differs from other kinds of software applications (Zhang et al., 2026). The behavior of artificial intelligence systems may be difficult to characterize using conventional metrics, requiring dedicated monitoring that reflects the semantic properties of model input and output. Retrieval-augmented generation systems convert incident material into dense vectors for similarity-based historical pattern matching to identify abnormal patterns (Zhang et al., 2026).

Monitoring for the GAI Profile includes looking for adversarial inputs trying to force the behavior of an LLM, outputs that appear to have harmful or inappropriate content, or resource usage patterns that deviate from the norm and may indicate abuse (NIST, 2024). The monitoring procedures should contain both automated alerts for known threats and human analysis of unusual patterns.

8.3 Evidence-Ready Operations

Mature platforms generate evidence of their state during normal operations without specific evidence-collection efforts, which helps demonstrate compliance. Research on automatic compliance benchmarks considers that evidence generation should ideally be integrated into the operational workflow (Sharma et al., 2024). The rapid provision of evidence artifacts by automated compliance tools reduces the need for manual documentation.

Documentation requirements of the NIST Artificial Intelligence Risk Management Framework serve governance and accountability (NIST, 2023). Organizations can develop and maintain evidence artifacts showing the artificial

intelligence system is operating within defined parameters, that risk mitigation controls are operating as intended, and that governance processes are effectively providing oversight. Platform evidence systems should include specific requirements for artificial intelligence in addition to cybersecurity controls.

Continuous compliance monitoring can automate evidence collection and gap identification, reducing manual labor for audit purposes and speeding up drift detection. For example, applying CIS Benchmark compliance monitoring to Kubernetes worker nodes was shown to lead to a 95 percent reduction in persistent benchmark violations compared to periodic compliance assessment (Sharma et al., 2024).

Component	Data Sources	Security Value	Evidence Artifacts
Kubernetes Audit	API server interactions	Authentication, authorization, resource changes	Access logs, privilege escalation detection
Node Telemetry	OS audit, file integrity, network flows	Host compromise detection, lateral movement	System integrity reports, network baselines
Container Runtime	Process execution, file access, connections	Container escape, anomalous behavior	Runtime violation records, behavioral analysis
Workload Observability	Application logs, traces, metrics	Application-layer attacks, compromise indicators	Service health records, anomaly detection
AI Workload Monitoring	Inference logs, resource usage, data access	Prompt injection, data extraction, model abuse	Usage patterns, access audits, policy compliance

Table 3: Components and Evidence Artifacts of Observability [Souppaya et al. (2017); Torkura et al. (2017); Sharma et al. (2024); NIST (2023)]

9. Secure Artificial Intelligence Workload Patterns

9.1 Model Training and Pipeline Security

AI workloads include several patterns that have different security properties and therefore require architectural consideration. For example, training pipelines ingest data, perform computation-intensive processing, and produce model artifacts. As the NIST Artificial Intelligence Risk Management Framework defines data quality and integrity as foundational to the trustworthiness of artificial intelligence systems (NIST, 2023), security controls for data protection, isolation of compute infrastructure, and the integrity of artifacts emitted by a model should be considered.

Enforcing the principle of least privilege through workload identity grants instead of shared credentials is useful in minimizing access. A zero trust architecture study found that access to datasets should be at a fine granularity where identity and access context are used for authorization (Phiayura & Teerakanok, 2023). Encrypt data in transit between storage and compute nodes, and on persistent storage. Use data lineage tracking to demonstrate compliance for regulated data and assist incident investigation.

According to container security best practices (Souppaya et al., 2017), the pipeline infrastructure should be isolated from production serving environments so that experimental workloads do not have access to production data or credentials. Namespaced isolation with restricted egress prevents data leaking from the training environment, while resource quotas prevent training jobs from exhausting cluster capacity for production services.

The integrity of model artifacts must be maintained from the completion of training through to deployment. Model provenance is a governance requirement for verifying the origin of deployed models in the Generative Artificial Intelligence Profile (NIST, 2024). Signed models allow for version control, where only authorized outputs from the training pipeline are deployed. Model registries should use access control to manage the service accounts used for pushing new model versions and pulling approved ones.

9.2 Model Serving and Inference Security

Real-time model serving exposes inference endpoints to untrusted input, unlike typical web services. Adversarial input manipulation has been identified as a major category of risk for deploying large language models (LLMs) in the NIST

Generative Artificial Intelligence Profile (NIST, 2024). Inference endpoints are subject to prompt injection attacks, in which adversarial input attempts to override system instructions, extract training data, or otherwise manipulate model behavior. Mitigations include application-layer input validation and output filtering, in addition to platform controls such as output monitoring and isolation.

Networking should restrict inference endpoints to legitimate data sources and services, and container security research shows that network isolation provides security beyond what is provided at the application layer (Sultan et al., 2019). Using egress filtering to prevent model-serving pods from opening connections to external resources not contained in the allow-list limits the damage that can result from a successful prompt injection attack exfiltrating data or accessing an internal service. Rate limiting and resource quotas prevent denial of service via resource exhaustion of expensive inference workloads.

Research on AI operations has identified a number of domain-specific operational practices for LLMs (Zhang et al., 2026), including monitoring model behavior, evaluating model responses, and identifying new failure modes. Tool-augmented approaches enable large language models to interface with monitoring systems that perform diagnostic operations and automatic remediation (Zhang et al., 2026). Such operational capabilities should be accompanied by appropriate telemetry, logging, and alerting capabilities in the platforms.

9.3 Data Protection and Compliance

Three major areas of concern in data protection for artificial intelligence workloads are compliance with data protection laws, protection of intellectual property, and maintenance of operational integrity. The NIST Artificial Intelligence Risk Management Framework identifies privacy as a cross-cutting issue (NIST, 2023). Regulated data require safeguards maintaining these requirements regardless of whether the data is handled in traditional or AI pipelines.

The Generative Artificial Intelligence Profile identifies potential data privacy risks with generative AI, including the risk of large language models memorizing their training data and unintentionally disclosing it regardless of data access management (NIST, 2024). The Profile encourages organizations to consider memorization capability when selecting models and to filter outputs from models operating over sensitive domains.

Separating sensitive data from non-production environments may help reduce exposure of sensitive data. A zero trust architecture review determined that data classification should drive access controls (He et al., 2022). Production data should not be copied to development and testing environments unless masked or anonymized, and model evaluation on production data should be performed in production-equivalent security environments rather than on developer workstations or shared test clusters.

10. Empirical Context and Industry Metrics

10.1 Container Security Challenge Distribution

Analysis of practitioner literature provided numerical context for the priorities of security architecture. A multi-vocal literature review of 321 internet artifacts related to Kubernetes found 15 challenges reported by practitioners (Shamim et al., 2022). Most artifacts (41.1%) reported issues with lack of security practices and tools. Next, attack surface reduction (39.2%), lack of diagnostics tools (37.7%), maintenance (33.9%), and learning curve (28.3%) challenges were reported in the artifacts. Security challenges were the most common in the distribution and appeared most frequently in every artifact class. Concerns related to networking and migration cost were present in 26.2 percent of artifacts (Shamim et al., 2022).

The difference between practitioner and research perspectives appears to be important, with only 7 out of 105 peer-reviewed works specifically addressing security (Shamim et al., 2022). For 9 of the 15 identified challenges, no papers addressing them were found. These challenges include attack surface reduction, maintenance challenges, learning curve, migration cost, storage, testing, cultural change, idempotency, and hardware compatibility. This result suggests that some practitioner security challenges have yet to be addressed by the academic community, and practice-oriented reference architectures could fill this gap.

10.2 Compliance Automation Effectiveness

To measure compliance automation empirically, policy enforcement for CIS Benchmark compliance on Kubernetes worker nodes was investigated by deploying a prototype system on different cloud configurations (Sharma et al., 2024).

The overall process latency was reduced by 99.8 percent. Manual remediation had a mean time-to-compliance of over 120 minutes, whereas automated enforcement had a mean time-to-compliance of less than two minutes.

Stability of compliance posture demonstrated continual enforcement improves compliance. While conducting continuous load testing deploying 500 pods per hour over 24 hours, the compliance posture remained stable at 99.5 percent with automated enforcement, but drifted to 78 percent without automatic enforcement (Sharma et al., 2024). Compared to manual auditing, there is a 95% reduction in benchmark failures. The security automation tool implements 52 of the 65 worker-node-relevant controls from CIS Benchmark version 1.8, using 35 remediation mechanisms and 17 admission control policies, achieving automation coverage of 80%. The additional latency caused by this performance overhead was negligible, with a 99th percentile pod startup latency increase of 142 milliseconds, and less than 1 percent of CPU and less than 50 megabytes of memory per node (Sharma et al., 2024).

10.3 Effectiveness of Container Security Mechanisms

Experimental studies of containers have been conducted to quantify the effectiveness of container security mechanisms in the context of defense-in-depth. For example, in a study of 223 container-targeting attacks, 56.8% were successful against default container configurations (Sultan et al., 2019). This finding shows that default configuration is insufficient, and explicit hardening to the level of the reference architecture is needed.

More than 35 percent of system call traffic was removed via seccomp profiles while introducing insignificant performance cost (Sultan et al., 2019). The security namespace, along with container-specific Linux Security Module profiles, allowed for per-container security policies to be applied with only a 0.7 percent increase in latency. These performance characteristics indicate that each of the security controls in the reference architecture can be implemented without significant impact.

10.4 Artificial Intelligence Operations Metrics

Research into AI operations has shown that new monitoring approaches are needed: a survey on large language model operational use found that it was difficult to characterize large language model behavior using standard metrics (Zhang et al., 2026). The evaluation sets include 7,184 multiple-choice questions, 1,736 multi-operation question-answering questions, and 130 natural language prompts for script generation for different operations and automation scenarios (Zhang et al., 2026).

These metrics suggest that monitoring AI workloads requires functionality beyond traditional application observability. The architectural patterns in this reference architecture address semantic monitoring of model behavior in addition to operational monitoring of applications.

11.1 Expert Evaluation Framework

The validation of the reference architecture is performed based on design science research methodology with criteria that can be grouped into three dimensions: security challenge coverage, framework alignment, and implementability. Security challenge coverage examines how well the framework addresses known concerns of practitioners, its alignment with security standards, and its implementability.

11.2 Security Challenge Coverage Assessment

Using multi-vocal literature review, the architecture provides a coverage assessment mapping architectural components to the set of 15 practitioner challenges identified by Shamim et al. (2022).

Lack of security practices and tools (41.1 percent of artifacts): This architecture provides policy-driven governance, layered isolation, and evidence-ready operations. It includes controls such as pod security standards, admission control, secrets management, and continuous compliance monitoring.

Attack surface reduction (39.2 percent of artifacts): The architecture supports attack surface reduction through defense-in-depth across trust boundaries, workload classification for privilege isolation, and network segmentation for limiting communication paths. AI-specific addressing for attack surface is included for accelerator nodes.

Lack of diagnostics tools (37.7 percent of artifacts): The security monitoring architecture provides visibility into Kubernetes audit signals, node signals, container runtime signals, and workload signals through artificial intelligence workload monitoring extensions.

Maintenance (33.9 percent of artifacts): maintenance in the implementation roadmap is addressed through rollout phases, operations maturity phases, and continuous compliance monitoring implementing drift detection.

Learning curve (28.3 percent of artifacts): the architecture provides progressive workload classification frameworks, isolation tier definitions, and phased implementation, which is easier to learn than ad-hoc security configuration.

11.3 Framework Alignment Assessment

The architecture aligns with established security frameworks to validate design decisions against standard security guidance.

NIST Application Container Security Guide alignment: the architecture addresses all five risk categories defined by Souppaya et al. (2017): image, registry, orchestrator, container, and host. Hardening the control plane addresses the orchestrator risk category. Hardening of worker nodes addresses risks at the host level, policy-based governance addresses risks at the container layer, and supply chain concerns address image and registry risks.

NIST Artificial Intelligence Risk Management Framework alignment: This architecture implements artificial intelligence asset classification based on data integrity, model provenance, and inference endpoint protection as specified by NIST (2023), as well as evidence-ready operations supporting governance and accountability.

Alignment with zero trust architecture: The architecture aligns with the three primary pillars of zero trust as outlined by Phiyayura and Teerakanok (2023): identity governance through workload identity, network segmentation through default-deny policies, and continuous verification via compliance monitoring.

Policy-based governance is achieved by mapping Pod Security Standards against CIS workload configuration controls and mapping CIS node configuration controls against worker node hardening. Compliance automation, found to be 80 percent feasible, validates the architecture (Sharma et al., 2024).

11.4 Implementability Assessment

These architectural components have been shown to be implementable, whereby compliance automation work reduced time to compliance by 99.8 percent, and persistent failures by 95 percent using architectural patterns included in the reference architecture (Sharma et al., 2024). Security controls did not add significant overhead to system performance, with less than 1 percent increase in CPU usage and a 142-millisecond increase in latency, making it operationally viable.

The phased implementation roadmap provides structured deployment guidance that addresses practitioner concerns about learning curve, migration cost, and introduces foundational controls before deploying operationally mature and advanced capabilities. This lowers the risk associated with implementation by building capabilities incrementally.

12. Limitations

12.1 Architectural Scope Limitations

While this reference architecture is focused on multi-tenant Kubernetes clusters, not all deployments will require the same level of control and security, such as single-tenant dedicated clusters, or edge deployments that may not have sufficient resources for monitoring and policy enforcement. Organizations should assess the architecture's applicability in their deployment contexts and adjust controls accordingly.

This architecture applies to standard Kubernetes distributions, but may differ with managed Kubernetes services, which restrict access to individual nodes or provide security features that supplement or replace architecture components depending on the service capabilities and the security requirements of organizations.

12.2 Threat Model Limitations

This architecture addresses many commonly used attack patterns; however, it does not address all attacks. Insider threats specifically resulting from adversaries gaining privileged access via insiders (e.g., privileged admins) require additional organizational controls such as segregation of duties, privileged access management, and anomaly monitoring to mitigate these threats.

Compromised continuous integration and continuous deployment pipelines may have access to deployment automation, enabling attackers to circumvent platform architectural controls. Supply chain security for build systems and deployment tools requires complementary organizational security controls beyond platform architectural controls.

Well-resourced advanced persistent threats may combine multiple attack vectors to circumvent individual controls. Defense-in-depth provides multiple barriers to an attack but cannot guarantee attack prevention. Organizations should implement detection and response capabilities alongside preventive ones to mitigate risk from advanced persistent threats.

12.3 Performance and Operational Overhead

Policy enforcement and monitoring may introduce additional latency, which is problematic for latency-sensitive applications. One study of compliance automation found admission control caused a 142 ms 99th percentile increase in pod startup latency (Sharma et al., 2024). Where practical compensating controls are in place, organizations with strict latency requirements should examine control settings to best accommodate performance-sensitive workloads.

Platform engineering teams would also require capacity to implement and operate the security controls, and organizations without strong security expertise may find the architectural controls challenging to operate. The implementation roadmap reduces this issue through phased implementation, but cannot address the fundamental need for expertise.

12.4 Validation Limitations

Expert evaluation, while an indication of adequacy, is not empirical validation. We suggest that future work perform empirical evaluation by deploying the architecture in a pilot and measuring security and operational outcomes through case studies.

This work focuses on synthesizing existing frameworks and research results, and its application to AI workloads, not on developing new security mechanisms. Organizations with specialized security requirements may need to extend architectural components beyond documented patterns.

13. Implementation Roadmap and Future Directions

13.1 Phased Implementation Approach

Multiple implementation stages should be defined to iteratively build security capabilities that align with the maturity and operational capacity of an organization. A study on Kubernetes implementation found organizations reported difficulty implementing multiple security capabilities in parallel (Shamim et al., 2022). Gradual implementation allows teams to gain operational experience with each capability before implementing further controls.

Foundation controls are key controls that should be implemented as soon as possible. These controls deliver integrated identity access, multi-factor authentication for all admin access, network segmentation with default-deny, baseline pod security, and centralized logging. The Application Container Security Guide recommends implementing foundation controls before deploying workloads in production (Souppaya et al., 2017). Organizations should validate operational capability with baseline controls prior to advancing to advanced capabilities.

The operational maturity phase is where upgrade processes, patching procedures, incident logging and remediation, and tenant onboarding procedures are established. One study that examined automated compliance with benchmarks found that compliance is an ongoing activity rather than a one-time configuration task (Sharma et al., 2024). Runbooks should document common operational scenarios with security considerations.

The advanced capabilities phase typically comprises automated evidence collection, continuous compliance scanning, and improved runtime detections. AI-specific controls, such as model artifact governance, data pipeline monitoring, and inference endpoint protection may be deployed together or sequentially, based on an organization's priorities and the maturity of the underlying platform.

13.2 Organizational Considerations

Technical architecture requires organizational support; a study of the challenges of transitioning to a zero trust architecture reported that executive leadership sponsorship, cross-functional collaboration, and change management were success factors (Phiayura & Teerakanok, 2023). Zero trust requires organizational change, not just technology.

Governance structures should include platform governance and AI risk governance. The NIST Artificial Intelligence Risk Management Framework identifies governance as one of the functions for managing AI risk (NIST, 2023). Existing change advisory and security review processes, along with artificial intelligence governance requirements for model risk assessment and incident reporting, should accommodate Kubernetes-specific technical considerations.

Skills development addresses the need for knowledge to operate and secure platforms. Barriers to securely adopting Kubernetes have included lack of skills (Shamim et al., 2022). Container security, Kubernetes administration, and artificial intelligence operations are specialties requiring training resources for professionals.

13.3 Future Directions

Container platform security continues to develop, as do threats and countermeasures. Zero trust architecture is a trend in enterprise security, including for container platforms. Acceleration of zero trust implementation across government and enterprise has been noted, and organizations have increasingly begun treating zero trust as a planned objective (He et al., 2022).

Research in AIOps has observed that LLMs are increasingly being integrated into operations, such as in security operations centers (Zhang et al., 2026). New applications, such as log analysis and incident response, and the automated application of LLMs, can generate both operational improvements and security risks. Platform architectures need to accommodate such integration patterns.

The NIST Generative Artificial Intelligence Profile acknowledges that governance frameworks for AI systems will continue to evolve as technology capabilities and risk understanding improve (NIST, 2024). Organizations are encouraged to stay aware of regulatory evolution and design platform architectures supporting evolving regulatory requirements. Teams with product-oriented platforms, with clear standards, measurable outcomes, and evidence to support their operations, will be best positioned to adapt.

Conclusion

Enterprise Kubernetes platforms have become a shared product with security-by-default designs and operational governance. AI workloads introduce new security considerations, with high-value assets, such as training datasets, model artifacts, and inference endpoints, complementing traditional application services as part of the attack surface. The reference architecture addresses: comprehensive threat modeling, multi-tenant isolation strategies based on workload sensitivity, defense in depth for the Kubernetes control plane and worker nodes, policy-based enforcement through admission control and Kubernetes Pod Security Standards, and observability architectures to support security operations and compliance evidence gathering. The architecture responds to gaps between practitioner security needs and research guidance. Analysis of multi-vocal literature found security issues dominated practitioner discussions, but research only addressed security in seven of 105 peer-reviewed publications. Empirical evaluation of compliance automation shows that policy-based enforcement reduces time-to-compliance and the rate of violations beyond the initial compliance state. These results substantiate the architecture's vision of automated policy enforcement and runtime compliance monitoring.

The design science research methodology ensures systematic derivation from knowledge bases. Expert review is employed to ensure compliance with established security frameworks. Limitations include scope, threat model boundaries, and that findings are validated by expert assessment rather than production implementations. Follow-up research could assess case studies of production implementations and their security impact. Next generation platforms may focus more on identity-based authorization aligned with zero trust principles, improved software and model supply chain integrity, and better governance of artificial intelligence workloads. As artificial intelligence services become widespread, security may focus on artificial intelligence governance, with architectures demonstrating integrated management of model, data, and cyber risks. Companies who treat container platforms as products, define expectations, control operating outcomes with metrics, and document their controls can modernize while maintaining stakeholder trust and resiliency.

Disclaimer:

The views expressed in this article are solely those of the author and do not represent the views of American Express.

References

1. Anya Sharma, Marcus Chen, and Elena Petrova (2026). Automating CIS Benchmark Compliance for Kubernetes Worker Nodes in Cloud-Native Environments, <https://www.researchgate.net/profile/Lanre-Akintayo/publication/399614353>
2. Burns, B., Grant, B., Oppenheimer, D., Brewer, E., & Wilkes, J. (2016). Borg, Omega, and Kubernetes: Lessons learned from three container-management systems over a decade. *ACM Queue*, 14(1), 70-93. <https://queue.acm.org/detail.cfm?id=2898444>
3. Islam, Shamim. S., Alexander Gibson, J., Morrison, P., and Rahman, A. (2022). Benefits, Challenges, and Research Topics: A Multi-vocal Literature Review of Kubernetes. arXiv, <https://arxiv.org/pdf/2211.07032>
4. László KOVÁCS, (2018). NATIONAL CYBER SECURITY AS THE CORNERSTONE OF NATIONAL SECURITY, *Land Forces Academy Review*, DOI: 10.2478/raft-2018-0013
4. LINGZHE ZHANG et al. (2026). A Survey of AIOps in the Era of Large Language Models, *ACM Computing Surveys*, Volume 58, Issue 2. <https://doi.org/10.1145/3746635>
5. NIST (2023). Artificial Intelligence Risk Management Framework (AI RMF 1.0), <https://doi.org/10.6028/NIST.AI.100-1>
6. NIST AI (2024). "Artificial Intelligence Risk Management Framework: Generative Artificial Intelligence Profile," <https://doi.org/10.6028/NIST.AI.600-1>
7. PHIAYURA. P, TEERAKANOK . S, (2023). A Comprehensive Framework for Migrating to Zero Trust Architecture, *IEEE Access*, 2023. <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=10052642>
8. Souppaya, M., Morello, J., & Scarfone, K. (2017). *Application Container Security Guide*, National Institute of Standards and Technology. <https://csrc.nist.gov/csrc/media/publications/sp/800-190/draft/documents/sp800-190-draft.pdf>
9. Sultan, S., Ahmad, I., & Dimitriou, T. (2019). Container security: Issues, challenges, and the road ahead. *IEEE Access*, doi: 10.1109/ACCESS.2019.2911732
10. Tak, B. C., Isci, C., Duri, S., Bila, N., Nadgowda, S., & Kompella, J. (2017). Understanding security implications of using containers in the cloud. *USENIX Annual Technical Conference*, <https://www.usenix.org/system/files/conference/atc17/atc17-tak-paper.pdf>
11. Tavallaee, M., Bagheri, E., Lu, W., & Ghorbani, A. A. (2009). A detailed analysis of the KDD CUP 99 data set. *IEEE Symposium on Computational Intelligence for Security and Defense Applications*, DOI:10.1109/CISDA.2009.5356528
12. Torkura, K. A., Sukmana, M. I. H., & Meinel, C. (2017). Integrating continuous security assessments in microservices and cloud-native applications. *Proceedings of the 10th IEEE/ACM International Conference on Utility and Cloud Computing*, <https://doi.org/10.1145/3147213.3147229>
13. Vayghan, L. A., Saied, M. A., Toeroe, M., & Khendek, F. (2021). A Kubernetes controller for managing the availability of elastic microservice-based stateful applications. *Journal of Systems and Software*, <https://doi.org/10.1016/j.jss.2021.110924>
14. Yuanhang He, Daochao Huang, Lei Chen , Yi Ni, and Xiangjie Ma. (2022). *A Survey on Zero Trust Architecture: Challenges and Future Trends*, Wiley, <https://doi.org/10.1155/2022/6476274>