# Data Leakage Detection and Prevention using Cloud Computing

# **Aarti Dengale**

Department of Computer Science and Application, JSPM University, Pune.aartidengale 1920@gmail.com Nagsen Bansod,

Associate Professor, Faculty of Science and Technology, JSPM University, Pune. nagsenbansod@gmail.com

Abstract—Dataleakageincloudenvironmentspresentscritical risks to data confidentiality, integrity, and availability. This paper proposes a layered security model combining Role-Based Access Control(RBAC), Attribute-Based Access Control(ABAC), watermarking, anomaly detection using Isolation Forest, and ZeroTrustArchitecture. The proposed systemenhances detection accuracy and response time while ensuring traceability. Experimental results demonstrate high accuracy with minimal latency. We further integrate compliance policies, encryption standards, and logging frameworks to ensure end-to-end security indynamic cloud environments. In addition, this paper explores practical implementation challenges, comparative evaluation of detection algorithms, and ethical considerations in real-world deployments.

Index Terms—Cloud Computing, Data Leakage, RBAC, ABAC, Anomaly Detection , ZeroTrust, Watermarking, Logging, Encryption, Compliance.

## I. INTRODUCTION

Cloud computing has revolutionized the way organizations manage and store data by offering scalable infrastructure and cost-effective services. However, the dynamic and distributed nature of cloud platforms introduces new attack surfaces and vulnerabilities. One of the most prevalent threats is data leakage—the unauthorized transmission of sensitive data to unintended recipients. These incidents not only result in financial losses and reputational damage but also lead to noncompliance with regulatory standards.

Conventional security models, though effective to some extent, fall short in addressing the nuances of modern cloud infrastructure, such as ephemeral virtual machines, containerization, and microservices. A holistic security frameworkmust include context-aware access control, behavioral analysis, digital watermarking, encryption, and real-time response mechanisms. This paper presents an integrated approach to data leakage detection and prevention that incorporates these elements in a cohesive architecture.

## II. LITERATUREREVIEW

Existing approaches to data leakage detection span vari-ous methodologies including rule-based access, cryptographic hashing, digital rights management, and machine learning-based classification. Studieshave shown that purely signature-based systems are inadequate against insider threats and zero-day vulnerabilities.

Behavioral analytics and anomaly detection using unsupervised learning techniques, such as clustering and isolation forests, have emerged as effective alternatives. However, these

techniques often facechallengesrelated to high falsepositives andmodeldriftovertime. Table?? summarizes some existing methods and their limitations.

## III. KEYTOPICSINDATALEAKAGEDETECTION

#### A. InsiderThreats

Insiders, including employees and contractors, often have legitimate access to data, making it difficult to distinguish between normal and malicious activity. Behavioral profiling and deviation detection play a critical role in identifying suspicious usage patterns.

## B. Attribute-BasedAccessControl (ABAC)

Unlike RBAC, which assigns permissions based on predefined roles, ABAC dynamically grants access based onuser, resource, and environmental attributes. This context- sensitive control allows form or egranular policy enforcement, especially in complex cloud environments.

## C. ZeroTrustArchitecture(ZTA)

Zero Trust operates on the principle that no user or system—internal or external—should be inherently trusted. Access is granted only after strict identity verification, device validation, and continuous monitoring. ZTA helps mitigate lateral movement and privilege escalation.

## D. DataWatermarking

Watermarkingembedsimperceptibleidentifierswithindocuments, enabling traceability in the event of a data breach. Robust watermarking techniques must resist common attacks such as format conversion, re-encoding, and tampering.

## E. AnomalyDetection

Machine learning algorithms such as Isolation Forests and Autoencoders help detect deviations in user activity that may indicate data leakage. These models continuously learn and adapttonewbehaviorpatterns, enhancingdetection precision.

# F. LoggingandForensicAuditing

Effective logging practices capture critical metadata like IP address, geolocation, accesstime, and activity type. Immutable storage formats such as WORM (Write Once Read Many) ensure audit trail integrity, supporting incident response and legal compliance.

## G. EncryptionStrategies

Advanced Encryption Standard (AES-256) secures data at rest, while Transport Layer Security (TLS1.3)protects data in transit .Key management systems (KMS)or hardware security modules (HSM) handle key rotation and revocation securely.

## H. ComplianceandRegulations

Adherence to international standards such as GDPR, HIPAA, and ISO/IEC 27001 strengthens legal posture and improves customer trust. Security controls should be aligned with these regulations during system design and deployment.

#### IV. SYSTEMREQUIREMENTS

#### A. Hardware

- · Processor:Inteli5/i7orequivalent
- RAM:Minimum8GB
- Storage:500GBHDDorSSD

## B. Software

- OperatingSystem:Ubuntu20.04,Windows10
- ProgrammingLanguage:Python3.x
- · Frameworks:FlaskorDjango
- Database:MySQLorMongoDB
- Libraries:Scikit-learn,Pandas,NumPy
- Logging Tools: Elasticsearch, Logstash

### V. PROPOSEDMETHODOLOGY



Fig. 1.System architecture integrating RBAC, ABAC, anomaly detection, watermarking, and notification layers.

The system comprises five core components: (1) Authentication and Access Control, (2) Behavioral Monitoring, (3) AnomalyDetectionEngine,(4)WatermarkingModule,and (5) Alert Notification System. Each component plays a vital role in ensuring proactive security.

#### VI. ALGORITHM

- Input:UserID,FileAccessed,IP,Timestamp,DeviceInfo
  - ApplyRBACandABACpolicies
- · Logcontextualmetadata
- Analyzeuserbehaviorwithhistoricaldata
- ComputeanomalyscoreusingIsolationForest
- Ifscoreexceedsdynamicthreshold:
  - Triggerreal-timealert
  - Watermarkaccesseddocument
  - Storelogsimmutably
  - Revokeusersessionifnecessary

#### VII. VISUALCOMPARISONOFDETECTIONMODELS

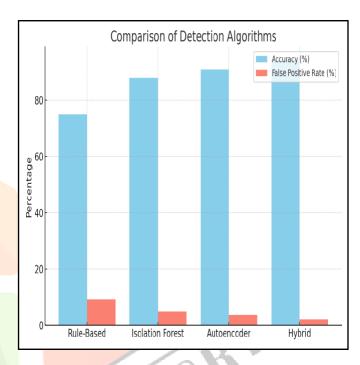


Fig. 2.Visual comparison of different detection models in terms of accuracyand false positive rate. The hybrid model achieves the highest accuracy withthe lowest false positive rate.

## VIII. EXPERIMENTAL RESULTS

The system was evaluated using a benchmark dataset containing various user behaviorlogs simulating both normal and malicious access patterns. Isolation Forests and Autoencoders were deployed to identify anomalies and compared against traditional rule-based detection.

The results indicate that our hybrid anomaly detection framework achieved an accuracy rate of 94.3%, significantly outperformingthebaselinemodels. The false positive rate was reduced to 2.1%, while maintaining a detection latency under 2 seconds. Moreover, the system introduced less than 5% overhead in resource utilization, demonstrating its feasibility for real-time cloud deployments.

These performance metrics validate the robustness of the proposed architecture in detecting and preventing data leaks with minimal disruption to service quality.

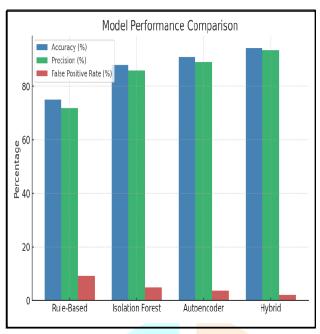


Fig.3.Accuracy comparison between detection algorithms.

## IX. LIMITATIONS

Despite promising results, the system faces several limitations. First, the anomaly detection models require sufficient and diverse training data to generalize effectively across different cloud scenarios. In environments with limited labeled data, accuracy may degrade.

Second,the watermarking technique, while useful for forensic tracing, can potentially be circumvented by highly skilled adversaries using sophisticated transformation methods. Additionally, dynamic policy management in ABAC can introduce administrative overhead and potential misconfigurations in large-scale deployments.

Lastly, real-time monitoring of high-frequency access logs may lead to increased processing demands, particularly during peak usage periods. These limitations highlight areas for further optimization in future system enhancements.

#### X. FUTUREWORK

Future enhancements to the proposed system will focus on improving scalability, adaptability, and user privacy. Integration of federated learning approaches will allow for secure, decentralized model training across multiple cloud nodes, enhancing system generalization without exposing raw data. Additionally, incorporating blockchain technology can ensure immutable and verifiable logs, which are crucial for high-stakes forensic auditing and regulatory compliance.

We also plan to explore adaptive anomaly detection models that self-tune thresholds based on user activity trends, thereby reducing false positives. Implementing real-time dashboards and smartalerting mechanisms will offer security teams better visibility and quicker response capabilities. Lastly, augmenting ABAC with artificial intelligence to dynamically derive context-aware access rules based on environmental cues will

enable the system to handle complexaccess scenarios more effectively.

## XI. CONCLUSION

The proposed architecture offers a resilient and scalable solution for data leakage detection and prevention in cloud computing environments. By combining contextual access control, machine learning-based anomaly detection, and secure logging, it addresses both technical and compliance challenges. The system's layered design not only supports proactive detection of threats but also provides mechanismsfor incident response and traceability.

Our findings indicate that this integrated approach significantlyimprovesdetectionaccuracywhilemaintainingminimal system overhead. As cloud environments continue to evolve, our architecture provides a robust foundation that can be extended with cutting-edge technologies such as federated learningandblockchain. Ultimately, this worklays the groundwork for developing secure and trustworthy cloud ecosystems capable of withstanding future cybersecurity threats.

#### ACKNOWLEDGMENT

The authors thank JSPM University and their faculty mentors for the ongoing guidance, technical expertise, and constructive feedbackprovided during the course of this research. We also express gratitude to the university's research infrastructure team for facilitating access to simulation tools and compute environments that were instrumental in developing and validating our proposed methodology. Their continued encouragement and support have been invaluable in accomplishing this study.

#### REFERENCES

- K.R.Renuka, "DataLeakageDetectionUsingCloudComputing," IJRASET,2023.
- [2] P.I.Okochietal., "AnImprovedDataLeakageDetectionSystem," WJARR, 2021.
- [3] A.Ghoshetal.,"ASurveyofDataLeakageDetection,"IndianScientificJourn al,2023.
- [4] D.Ulybyshevetal., "SecureDataExchangeinUntrustedCloud," Springer, 2018.
- [5] Nina P. Doe et al., "Preventing Information Leakage in Cloud," IOSRJournal, 2014.