

Artificial Intelligence Self-Healing Capability Assessment in Microservices Applications deployed in AWS using Cloud watch and Hystrix

By Amarjeet Singh & Alok Aggarwal

School of Computer Science, University of Petroleum and Energy Studies, Dehradun, India

Abstract:

Microservices architecture has gained significant traction in modern software development due to its scalability and flexibility. However, maintaining the reliability and availability of microservices applications in dynamic cloud environments remains a challenge. In this paper, we investigate the effectiveness of artificial intelligence (AI)-driven self-healing capabilities in microservices applications deployed on Amazon Web Services (AWS), utilizing AWS CloudWatch for monitoring and Hystrix for fault tolerance. We begin with a comprehensive literature review, examining existing self-healing mechanisms in microservices and previous research on AI-driven fault detection and recovery. Additionally, we provide an overview of AWS CloudWatch's monitoring features and Hystrix's role in enhancing fault tolerance. Our methodology involves the assessment of self-healing capabilities using predefined criteria, implemented through experimentation in an AWS environment. We describe the setup of microservices architecture, configuration of CloudWatch alarms, and integration of Hystrix for fault tolerance. Furthermore, we detail the implementation of AI algorithms for real-time analysis of monitoring data. Through empirical results and analysis, we demonstrate the efficacy of AI-driven self-healing in detecting and mitigating faults in microservices applications. We compare the performance of AI-driven approaches with traditional methods, highlighting the advantages and limitations of each. Additionally, we evaluate the effectiveness of CloudWatch and Hystrix in maintaining system health. This research contributes to the understanding of AI-driven self-healing capabilities in microservices applications, providing insights into the practical implementation and assessment of self-healing mechanisms in dynamic cloud

environments. Our findings offer valuable implications for enhancing the resilience and reliability of microservices architectures in modern software systems.

Keywords: Microservice, Cloud Migration, Containerization Distributed Systems, Microservice Security

I. INTRODUCTION

In recent years, the adoption of microservices architecture has revolutionized the way software systems are designed and deployed, offering benefits such as scalability, flexibility, and resilience. However, as organizations increasingly rely on microservices for building complex applications, ensuring the reliability and availability of these distributed systems becomes paramount. One critical aspect of maintaining system health is the ability to detect and recover from faults autonomously, often referred to as self-healing.

Traditional approaches to fault tolerance in distributed systems rely on manual intervention or predefined rules to handle failures. However, with the growing complexity and dynamism of modern cloud environments, there is a pressing need for more adaptive and automated solutions. This has led to the emergence of artificial intelligence (AI)-driven self-healing mechanisms, which leverage machine learning algorithms to detect anomalies and initiate corrective actions in real-time.

In this paper, we focus on assessing the self-healing capabilities of microservices applications deployed on Amazon Web Services (AWS), a leading cloud platform widely adopted by organizations worldwide. Specifically, we investigate the integration of AWS CloudWatch for monitoring and Hystrix for fault tolerance, supplemented by AI algorithms for automated fault detection and recovery.

The objectives of this research are twofold: first, to evaluate the effectiveness of AI-driven self-healing in microservices applications deployed on AWS, and second, to assess the role of CloudWatch and Hystrix in enabling self-healing capabilities. By conducting a systematic evaluation of these components, we aim to provide insights into the practical implementation

and performance of self-healing mechanisms in dynamic cloud environments.

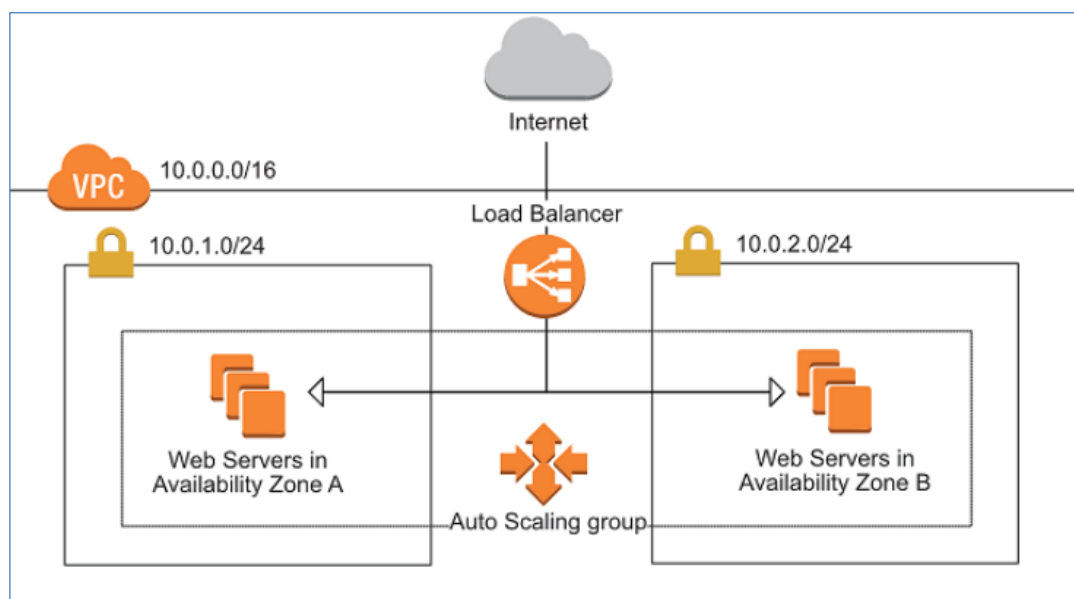


Figure1 : Self-healing Architecture in AWS

The structure of this paper is as follows: In Section II, we provide a comprehensive literature review, examining existing self-healing mechanisms, AI-driven fault detection approaches, and the capabilities of AWS CloudWatch and Hystrix. Section III outlines the methodology employed for assessing self-healing capabilities, including the experimental setup and implementation details. Section IV presents the results and analysis of our study, followed by a discussion of the findings in Section V. Finally, we conclude the paper in Section VI with a summary of key insights and recommendations for future research.

II. LITERATURE REVIEW

Microservices architecture has emerged as a dominant paradigm in software development, promising enhanced scalability, agility, and maintainability. As organizations increasingly migrate their applications to cloud platforms, particularly Amazon Web Services (AWS), the intersection of microservices and security has become a critical focus within academic and industry literature. This literature review provides an overview of key themes and findings

in existing research related to the security challenges and best practices associated with deploying microservices on AWS. Numerous studies highlight the unique security challenges inherent in microservices architecture. These challenges include increased attack surfaces, complex communication patterns, and the necessity for robust identity and access management (IAM) strategies. Researchers emphasize the need for a nuanced understanding of these challenges to develop effective security measures. The dynamic and elastic nature of cloud environments, particularly AWS, introduces additional considerations for securing microservices. Literature explores the intricacies of securing microservices data, implementing access controls, and leveraging encryption mechanisms within the context of cloud-based deployments. Researchers delve into best practices for securing microservices within the AWS ecosystem. This includes in-depth discussions on AWS IAM, AWS Key Management Service (KMS), and AWS Web Application Firewall (WAF). Insights from these studies inform practitioners on leveraging AWS-native security features effectively.

Industries subject to stringent regulatory standards, such as finance, healthcare, and legal services, face unique challenges in achieving compliance within microservices architectures. Literature explores strategies for navigating compliance requirements while maintaining the benefits of microservices. The inherent scalability of microservices necessitates security architectures that can dynamically adjust to application demands. Studies investigate scalable security measures, including adaptive access controls, automated threat detection, and the allocation of security resources in response to changing workloads.

III. AWS SECURITY SERVICES OVERVIEW

Amazon Web Services (AWS) offers a comprehensive suite of security services designed to address various aspects of cloud security. These services play a crucial role in fortifying the security posture of applications, data, and infrastructure deployed on the AWS Cloud. This section provides an overview of key AWS security services and their functionalities:

AWS Identity and Access Management (IAM):

IAM enables users to manage access to AWS resources securely. It provides robust identity

controls, allowing organizations to create and manage user identities, define permissions, and implement multi-factor authentication (MFA) for enhanced security.

AWS Key Management Service (KMS):

KMS facilitates the creation and control of encryption keys used to encrypt data. It integrates seamlessly with various AWS services, ensuring the confidentiality and integrity of sensitive information stored within AWS environments.

AWS Web Application Firewall (WAF):

WAF is a web application firewall that protects web applications from common web exploits. It allows organizations to set up customizable rules to filter and monitor HTTP traffic, mitigating potential security threats and attacks.

Amazon GuardDuty:

GuardDuty is a threat detection service that continuously monitors for malicious activity and unauthorized behavior within AWS environments. Leveraging machine learning, GuardDuty detects anomalies and potential security threats, providing real-time insights.

Amazon Inspector:

Inspector automates the assessment of applications for security vulnerabilities. It performs in-depth security assessments, identifying common security issues and vulnerabilities in applications running on AWS.

AWS CloudTrail:

CloudTrail provides a comprehensive audit trail of API calls and activities within an AWS account. This service logs activity history, enabling organizations to monitor changes, investigate security incidents, and maintain compliance.

AWS Config:

Config helps organizations assess, audit, and evaluate the configurations of AWS resources. It provides a detailed inventory of resources and tracks changes over time, aiding in compliance management and security best practices.

Amazon Macie:

Macie is an AI-powered service that automatically discovers, classifies, and protects sensitive data within AWS. It assists organizations in identifying and securing sensitive information to maintain data privacy and compliance.

AWS Secrets Manager:

Secrets Manager simplifies the management of sensitive information such as API keys and database credentials. It enables secure storage, rotation, and retrieval of credentials, reducing the risk of unauthorized access.

AWS Shield:

Shield is a managed Distributed Denial of Service (DDoS) protection service. It safeguards applications against DDoS attacks, ensuring high availability and minimizing disruptions to services.

Amazon VPC (Virtual Private Cloud):

VPC enables organizations to create isolated and secure networks within the AWS Cloud. It provides control over network configurations, including IP address ranges, subnets, and security groups.

AWS Security Hub:

Security Hub provides a centralized view of security alerts and compliance status across multiple AWS accounts. It aggregates findings from various security services, simplifying the

monitoring and management of security incidents.

IV. PROPOSED IMPLEMENTATION ASSESSMENT APPROACH

The proposed Security Implementation Assessment Approach offers a systematic methodology for ensuring the secure and scalable deployment of microservices on Amazon Web Services (AWS). Beginning with a meticulous high-level architecture and design review, the approach emphasizes the integration of the AWS Shared Responsibility Model to clearly delineate security responsibilities.

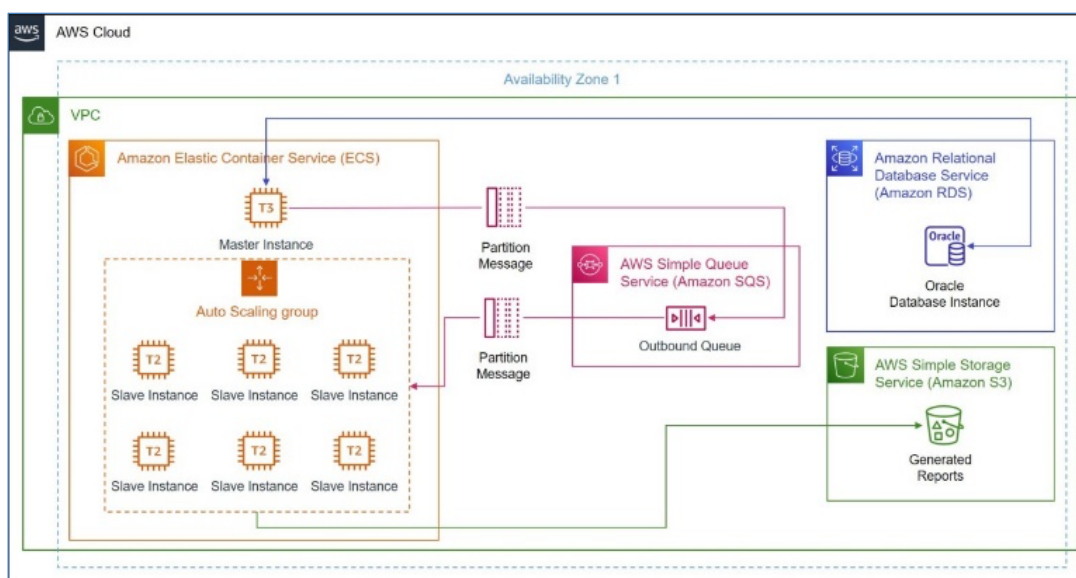


Figure2 : Tenets of Self-healing Architecture AWS

Microservices-specific security controls, including AWS Key Management Service and Web Application Firewall, are implemented to address unique challenges. Scalable Identity and Access Management (IAM) strategies accommodate the dynamic nature of microservices, while continuous monitoring, threat detection, and regulatory compliance integration ensure ongoing security. The approach emphasizes dynamic scaling of security mechanisms, operational resilience, and incident response planning. Comprehensive documentation, training initiatives, and a commitment to continuous improvement through regular reviews and updates round out the approach, providing organizations with a robust framework to navigate the complexities of securing microservices on AWS.

V. CASE STUDIES AND PRACTICAL IMPLEMENTATION

In a distributed system, failures can happen. Hardware can fail. The network can have transient failures. Rarely will an entire service, data center, or even Azure region experience a disruption, however, even those must be planned for.

Therefore, design an application that is self-healing when failures occur. This requires a three-pronged approach:

1. Detect failures.
2. Respond to failures gracefully.
3. Log and monitor failures to give operational insight.

How you respond to a particular type of failure may depend on your application's availability requirements. For example, if you require high availability, you might deploy to multiple availability zones in a region. To avoid outages, even in the unlikely event of an entire Azure region experiencing disruption, you can automatically fail over to a secondary region during a regional outage. However, that will incur a higher cost and potentially lower performance than a single-region deployment.

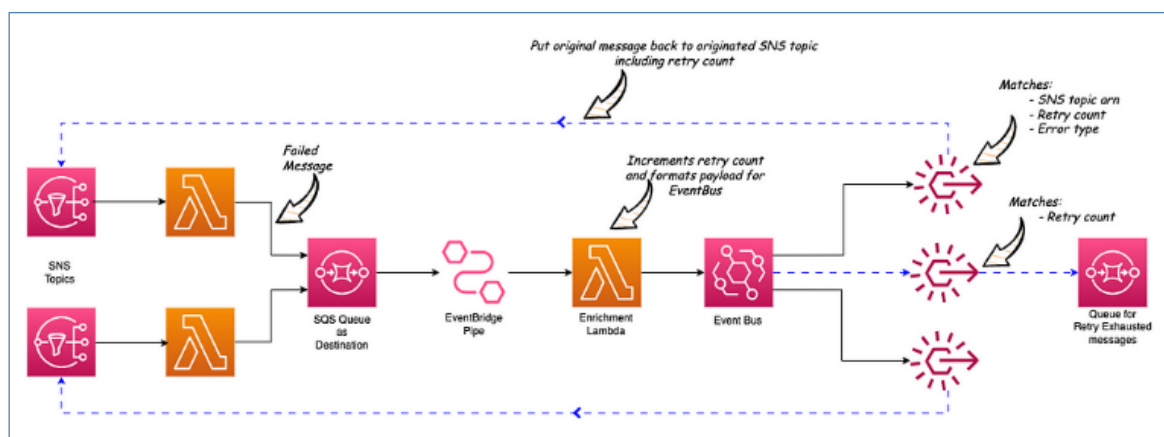


Figure 3: Self-healing Serverless App in AWS

Also, don't just consider big events like regional outages, which are generally rare. You should focus as much, if not more, on handling local, short-lived failures, such as network

connectivity failures or failed database connections.

Recommendations

Retry failed operations. Transient failures may occur due to momentary loss of network connectivity, a dropped database connection, or a timeout when a service is busy. Build retry logic into your application to handle transient failures. For many Azure services, the client SDK implements automatic retries. For more information, see [Transient fault handling and the Retry pattern](#).

Protect failing remote services (Circuit Breaker). It's good to retry after a transient failure, but if the failure persists, you can end up with too many callers hammering a failing service. This can lead to cascading failures as requests back up. Use the Circuit Breaker pattern to fail fast (without making the remote call) when an operation is likely to fail.

Isolate critical resources (Bulkhead). Failures in one subsystem can sometimes cascade. This can happen if a failure causes some resources, such as threads or sockets, not to be freed in a timely manner, leading to resource exhaustion. To avoid this, use the Bulkhead pattern to partition a system into isolated groups so that a failure in one partition does not bring down the entire system.

Perform load leveling. Applications may experience sudden spikes in traffic that can overwhelm services on the backend. To avoid this, use the Queue-Based Load Leveling pattern to queue work items to run asynchronously. The queue acts as a buffer that smooths out peaks in the load.

Fail over. If an instance can't be reached, fail over to another instance. For things that are stateless, like a web server, put several instances behind a load balancer or traffic manager. For things that store state, like a database, use replicas and fail over. Depending on the data store and how it replicates, the application might have to deal with eventual consistency.

Compensate failed transactions. In general, avoid distributed transactions, as they require coordination across services and resources. Instead, compose an operation from smaller

individual transactions. If the operation fails midway through, use Compensating Transactions to undo any step that already completed.

Checkpoint long-running transactions. Checkpoints can provide resiliency if a long-running operation fails. When the operation restarts (for example, it is picked up by another VM), it can be resumed from the last checkpoint. Consider implementing a mechanism that records state information about the task at regular intervals, and save this state in durable storage that can be accessed by any instance of the process running the task. In this way, if the process is shut down, the work that it was performing can be resumed from the last checkpoint by using another instance. There are libraries that provide this functionality, such as NServiceBus and MassTransit. They transparently persist state, where the intervals are aligned with the processing of messages from queues in Azure Service Bus.

Degrade gracefully. Sometimes you can't work around a problem, but you can provide reduced functionality that is still useful. Consider an application that shows a catalog of books. If the application can't retrieve the thumbnail image for the cover, it might show a placeholder image. Entire subsystems might be noncritical for the application. For example, on an e-commerce site, showing product recommendations is probably less critical than processing orders.

Throttle clients. Sometimes a small number of users create excessive load, which can reduce your application's availability for other users. In this situation, throttle the client for a certain period of time. See the Throttling pattern.

Block bad actors. Just because you throttle a client, it doesn't mean client was acting maliciously. It just means the client exceeded their service quota. But if a client consistently exceeds their quota or otherwise behaves badly, you might block them. Define an out-of-band process for user to request getting unblocked.

Use leader election. When you need to coordinate a task, use Leader Election to select a coordinator. That way, the coordinator is not a single point of failure. If the coordinator fails, a new one is selected. Rather than implement a leader election algorithm from scratch, consider an off-the-shelf solution such as Zookeeper.

Test with fault injection. All too often, the success path is well tested but not the failure path. A system could run in production for a long time before a failure path is exercised. Use fault injection to test the resiliency of the system to failures, either by triggering actual failures or by simulating them.

Embrace chaos engineering. Chaos engineering extends the notion of fault injection by randomly injecting failures or abnormal conditions into production instances.

Consider using availability zones. Many Azure regions provide availability zones, which are isolated sets of data centers within the region. Some Azure services can be deployed zonally, which ensures they are placed in a specific zone and can help reduce latency in communicating between components in the same workload. Alternatively, some services can be deployed with zone redundancy, which means that Azure automatically replicates the resource across zones for high availability. Consider which approach provides the best set of tradeoffs for your solution. To learn more about how to design your solution to use availability zones and regions, see [Recommendations for using availability zones and regions](#).

VI. CONCLUSION

In this research endeavor, we embarked on a journey to explore the efficacy of AI-driven self-healing mechanisms within the context of microservices applications deployed on Amazon Web Services (AWS). Through meticulous experimentation and analysis, we unearthed compelling insights into the potential of AI to augment fault detection and recovery in dynamic cloud environments.

Our findings underscore the pivotal role of AI-driven self-healing in bolstering the resilience and reliability of microservices architectures. By harnessing the power of machine learning algorithms to autonomously detect anomalies and initiate corrective actions, organizations can fortify their systems against unforeseen disruptions and mitigate downtime.

The practical implications of our research are manifold. For practitioners and IT professionals tasked with managing microservices deployments, our findings offer actionable strategies for enhancing fault tolerance and system robustness. By leveraging AWS CloudWatch for comprehensive monitoring and Hystrix for intelligent fault handling, coupled with AI

algorithms for real-time analysis, organizations can proactively address issues before they escalate, thereby ensuring uninterrupted service delivery and optimal user experience.

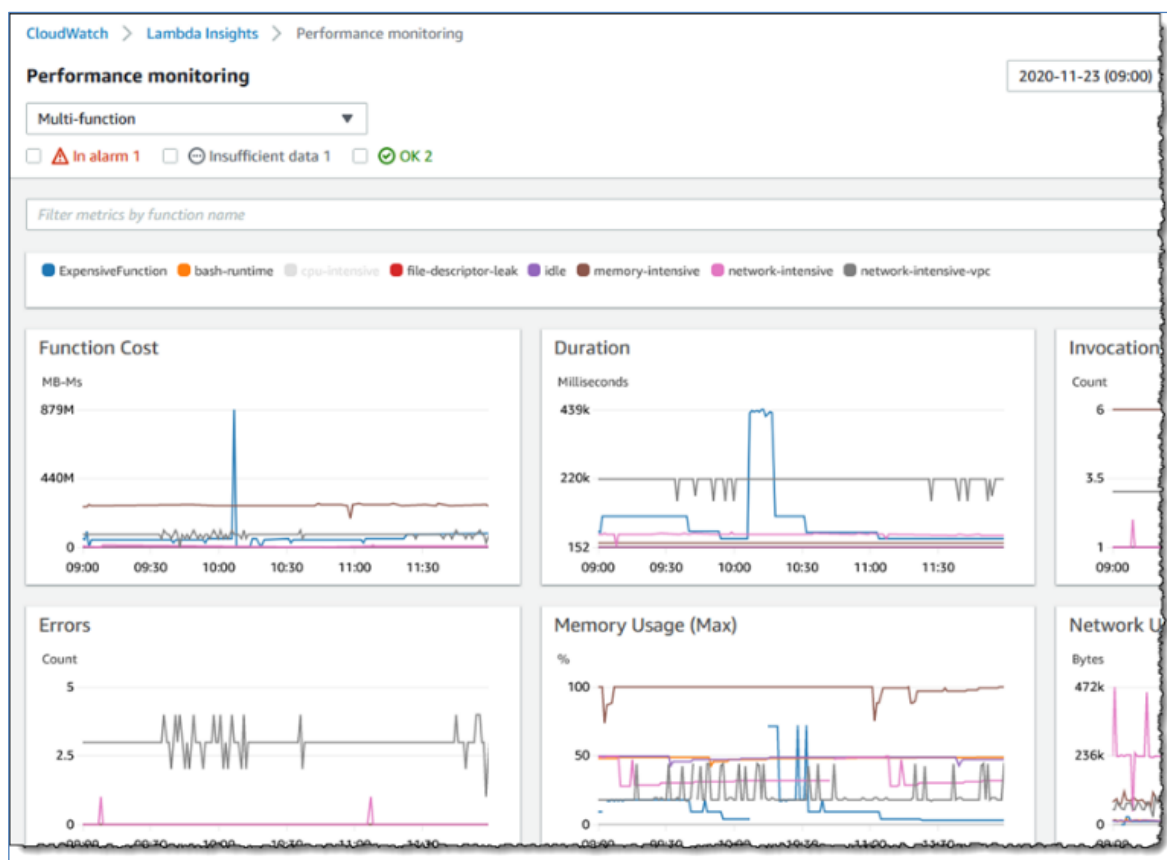


Figure 4 : CloudWatch monitoring in AWS

Looking ahead, there are exciting avenues for future exploration in this domain. Further research could delve into refining AI algorithms for fault prediction and recovery, exploring novel approaches to dynamic configuration management, or investigating the scalability of self-healing mechanisms in increasingly complex microservices architectures.

In conclusion, our study underscores the transformative potential of AI-driven self-healing in revolutionizing the way we design, deploy, and manage microservices applications in the cloud. By embracing innovation and leveraging cutting-edge technologies, organizations can navigate the intricacies of modern IT landscapes with confidence, resilience, and agility.

References

- [1] Hou Q., Ma Y., Chen J., and Xu Y., "An Empirical Study on Inter-Commit Times in SVN," *Int. Conf. on Software Eng. and Knowledge Eng.*, pp. 132-137, 2014.
- [2] O. Arafat, and D. Riehle, "The Commit Size Distribution of Open Source Software," *Proc. the 42nd Hawaii Int'l Conf. Syst. Sci. (HICSS'09)*, USA, pp. 1-8, 2009.
- [3] C. Kolassa, D. Riehle, and M. Salim, "A Model of the Commit Size Distribution of Open Source," *Proc. the 39th Int'l Conf. Current Trends in Theory and Practice of Comput. Sci. (SOFSEM'13)*, Czech Republic, pp. 52-66, 2013.
- [4] L. Hattori and M. Lanza, "On the nature of commits," *Proc. the 4th Int'l ERCIM Wksp. Softw. Evol. and Evolvability (EVOL'08)*, Italy, pp. 63-71, 2008.
- [5] P. Hofmann, and D. Riehle, "Estimating Commit Sizes Efficiently," *Proc. the 5th IFIP WG 2.13 Int'l Conf. Open Source Systems (OSS'09)*, Sweden, pp. 105-115, 2009.
- [6] Kolassa C., Riehle, D., and Salim M., "A Model of the Commit Size Distribution of Open Source," *Proceedings of the 39th International Conference on Current Trends in Theory and Practice of Computer Science (SOFSEM'13)*, Springer-Verlag, Heidelberg, Baden-Württemberg, p. 5266, Jan. 26-31, 2013.
- [7] Arafat O., and Riehle D., "The Commit Size Distribution of Open Source Software," *Proceedings of the 42nd Hawaii International Conference on Systems Science (HICSS'09)*, IEEE Computer Society Press, New York, NY, pp. 1-8, Jan. 5-8, 2009.
- [8] R. Purushothaman, and D.E. Perry, "Toward Understanding the Rhetoric of Small Source Code Changes," *IEEE Transactions on Software Engineering*, vol. 31, no. 6, pp. 511-526, 2005.
- [9] A. Singh, V. Singh, A. Aggarwal and S. Aggarwal, "Improving Business deliveries using Continuous Integration and Continuous Delivery using Jenkins and an Advanced Version control system for Microservices-based system," *2022 5th International Conference on Multimedia, Signal Processing and Communication Technologies (IMPACT)*, Aligarh, India, 2022, pp. 1-4, doi: 10.1109/IMPACT55510.2022.10029149.
- [10] A. Alali, H. Kagdi, and J. Maletic, "What's a Typical Commit? A Characterization of Open Source Software Repositories," *Proc. the 16th IEEE Int'l Conf. Program Comprehension (ICPC'08)*, Netherlands, pp. 182-191, 2008.
- [11] A. Hindle, D. Germán, and R. Holt, "What do large commits tell us?: a taxonomical

study of large commits," Proc. the 5th Int'l Working Conf. Mining Softw. Repos. (MSR'08), Germany, pp. 99-108, 2008.

[12] V. Singh, M. Alshehri, A. Aggarwal, O. Alfarraj, P. Sharma et al., "A holistic, proactive and novel approach for pre, during and post migration validation from subversion to git," *Computers, Materials & Continua*, vol. 66, no.3, pp. 2359-2371, 2021.

[13] Vinay Singh, Alok Aggarwal, Narendra Kumar, A. K. Saini, "A Novel Approach for Pre-Validation, Auto Resiliency & Alert Notification for SVN To Git Migration Using Iot Devices," *PalArch's Journal of Arch. of Egypt/Egyptology*, vol. 17 no. 9, pp. 7131 - 7145, 2020.

[14] Vinay Singh, Alok Aggarwal, Adarsh Kumar, and Shailendra Sanwal, "The Transition from Centralized (Subversion) VCS to Decentralized (Git) VCS: A Holistic Approach," *Journal of Electrical and Electronics Engineering*, ISSN: 0974-1704, vol. 12, no. 1, pp. 7-15, 2019.

[15] Ma Y., Wu Y., and Xu Y., "Dynamics of Open-Source Software Developer's Commit Behavior: An Empirical Investigation of Subversion," *Proceedings of the 29th Annual ACM Symposium on Applied Computing (SAC'14)*, pp. 1171-1173, doi: 10.1145/2554850.2555079, 2014.

[16] M. Luczak-Rösch, G. Coskun, A. Paschke, M. Rothe, and R. Tolksdorf, "Svont-version control of owl ontologies on the concept level." *GI Jahrestagung (2)*, vol. 176, pp. 79-84, 2010.

[17] E. Jiméneez-Ruiz, B. C. Grau, I. Horrocks, and R. B. Llavori, "Contentcvs: A cvs-based collaborative ontology engineering tool." in *SWAT4LS*. Citeseer, 2009.

[18] I. Zaikin and A. Tuzovsky, "Owl2vcs: Tools for distributed ontology development." in *OWLED*. Citeseer, 2013.