

Autonomous Incident Remediation: A Closed-Loop RAG Framework for Real-Time Root Cause Analysis and Knowledge Synthesis in Distributed Cloud Systems

Chetan Sasidhar Ravi¹ and Rohit Reddy Patlolla²

Integration Subject Matter Expert Zurich American Insurance Company, Schaumburg, IL USA

Integration Engineer, REI, Seattle, WA USA

chetan.ravi87@gmail.com and rohitredy@gmail.com

Abstract: *As cloud-native architectures grow in complexity, the “Mean Time to Recovery” (MTTR) is increasingly bottlenecked by the human capacity to parse massive volumes of unstructured log telemetry. This paper introduces an autonomous framework for Log-to-Lesson (L2L) synthesis, leveraging Retrieval-Augmented Generation (RAG) and Large Language Models (LLMs) to automate the lifecycle of incident response. Unlike traditional rule-based alerting, our proposed architecture implements a proactive “Consult-Research-Synthesize” (CRS) loop.*

Upon detecting a system anomaly, the framework first performs a semantic similarity search across a Cloud-Based Vector Store to identify historical Root Cause Analysis (RCA) reports or Knowledge Base (KB) articles. If a matching resolution is absent, the system triggers an “Agentic Research” phase: an LLM-based agent employing a ReAct (Reason+Act) loop parses raw log dumps, correlates error signatures with system metadata, and generates a novel, structured RCA report.

We demonstrate the efficacy of this framework through a 12-month deployment on a multi-cloud Kubernetes environment spanning GKE, EKS, and AKS clusters. Our results show that the system achieves a 92% accuracy rate in identifying recurring issues and reduces manual investigation time by 65%. Furthermore, the framework automatically populates a long-term knowledge repository in cloud storage, codifying “tribal knowledge” into structured digital assets. This study provides a blueprint for the transition from reactive monitoring to Cognitive Self-Healing, offering a scalable solution for Site Reliability Engineering (SRE) in the era of autonomous operations..

Keywords: Retrieval-Augmented Generation (RAG), Large Language Models (LLMs), Root Cause Analysis (RCA), AIOps, Cognitive Observability, Self-Healing Systems, Log Telemetry, Service Reliability Engineering, Vector Database, Agentic AI, Multi-Cloud Kubernetes