# AI-DRIVEN ROOT CAUSE ANALYSIS FRAMEWORK FOR DISTRIBUTED MICROSERVICES ARCHITECTURES

[1]Awodele S. O
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
awdeles@babcock.edu.ng

[2]Faruna, J. O
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
faruna0100@pg.babcock.edu.ng

[3]Mustapha M. M
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
mustapha0219@pg.babcock.edu.ng

[4]Ojuawo O. O
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
ojuawo0687@pg.babcock.edu.ng

[5]Olorunyomi O. B
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
Olorunyomi0052@pg.babcock.edu.ng

[6]Chukwulobe I
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
chukwulobe0408@pg.babcock.edu.ng

[7]Fayemi T. A
Department of Computer Science, Babcock University, Ilishan-Remo, Ogun State, Nigeria
fayemi0197@pg.babcock.edu.ng

**ABSTRACT**

Multi-layered IT systems are characterized by persistent system downtime, which creates significant operational burdens due to high costs and prolonged troubleshooting processes. Traditional root cause analysis (RCA) techniques are largely ineffective because they are predominantly reactive and manual, making them unsuitable for the scale and complex interdependencies of modern microservices architectures. This study addresses the critical "observability–complexity gap," a condition in which contemporary monitoring tools generate vast amounts of correlated data but fail to provide true causal resolution. To address this challenge, the study introduces RCASage, a novel AI-augmented inference engine designed for autonomous fault identification. The core innovation of RCASage lies in its hybrid architecture, which comprises a three-stage pipeline: (1) multi-modal telemetry data ingestion and dynamic dependency graph construction; (2) unsupervised anomaly detection using LSTM autoencoders combined with Natural Language Processing (NLP) classifiers; and (3) an Autonomous Inference Engine (AIE). This engine uniquely integrates Graph Neural Networks (GNNs) with the Neural Granger Causal Discovery algorithm to distinguish true root causes from downstream symptomatic effects. By shifting RCA from symptom correlation to causal inference, RCASage surpasses existing approaches such as beta-binomial inference and event-graph-based systems (e.g., GROOT). Furthermore, it bridges the gap between development and operations by incorporating Just-in-Time (JIT) defect prediction into the CI/CD pipeline. Empirical evaluation demonstrates that the proposed AI-augmented framework reduces Mean Time to Resolution (MTTR) by over 90% compared to traditional manual approaches. In addition, RCASage advances the RCA paradigm by embedding Explainable AI (XAI) principles, providing transparent causal explanations alongside visual dependency graphs. This positions RCASage as an intelligent, evidence-based, and autonomous solution for root cause analysis in contemporary digital infrastructures.

**Keywords:** Root Cause Analysis (RCA), RCASage, AI-Driven, AI-Augmented, Causal Inference, Distributed Systems, Graph Neural Networks, Graph-based Deep Learning, Anomaly Detection, Software Defect Prediction, Continuous Integration/Continuous Deployment (CI/CD), Microservices Architecture, Explainable AI (XAI), CI/CD Proactive Monitoring.

## 1    INTRODUCTION

Microservices architecture has become a de facto standard in the development of large-scale, resilient, scalable, and fast-changing distributed systems [1], [2]. Microservices architecture disaggregates applications into sets of bite-sized, independently deployable services that have many advantages over monolithic systems, including greater scalability, modularity, agility, and flexible continuous delivery [3], [4], [5]. Nevertheless, this complete architectural transformation dramatically raises the level of operational complexity [3]. With this kind of dynamism, individual user requests traverse a plethora of loosely coupled services that result in a high likelihood of possible points of failure, service communication, data integrity, and distributed monitoring systems, among others [1], [2], [6].

Root Cause Analysis (RCA) is believed to be one of the most difficult and time-consuming operations in modern IT infrastructures [6], [7], [8]. The problem of analytically tracing the system behaviour is exacerbated by the interconnectivity of the microservices, where a failure in one microservice can spread very fast, causing other microservices to fail, which obscures the original cause of the failure and leads to the overall failure of the distributed system [1], [7]. This dynamic gives rise to the "Observability-Complexity" paradox: although modern monitoring tools can produce large volumes of telemetry data (logs, metrics, traces) and provide a solid correlation, they may be inaccurate in inferring the root cause [6]. Therefore, the conventional RCA techniques, which rely on the inspection of these logs manually, the use of fixed rule-based correlation, and ad-hoc troubleshooting, are inefficient and are characterized by slowness, error, and, in general, are inefficient in scale, speed, and dependency, resulting in lengthy Mean Time to Resolve (MTTR) [6], [7], [9]. Thus, the following set of crucial issues emerges:

- Propagated Misdiagnosis: Monitoring tools tend to identify the most critical symptom of failure as the source of anomaly, but do not trace its origin to a specific component that may or may not be critical or unusual.
- Analytical Overload: The huge volume of available telemetry data overwhelms human operators with long diagnostic times and unacceptable Mean Time to Resolve (MTTR).
- Brittle and Noisy Operations: Rule-based AIOps systems are also a part of the issue by being inflexible in their threshold-based notifications and being unaware of subtle yet significant "gray failures".

In order to bridge this gap successfully, this paper explores how Artificial Intelligence (AI) and Machine Learning (ML) can be used to automate and enhance RCA [10], [11], [12], [13]. The future of AI-driven frameworks depends on their ability to transform RCA from a reactive, human-centered process into a proactive, data-centered science. This is done by using methods like causal inference, graph-based deep learning, and multi-modal anomaly detection to be able to autonomously reason over complex chains of failure, trying not only to report the symptoms, but also discover the real root cause [10], [11]. The other critical aspect is that these automated methods employ the principles of Explainable AI (XAI) that ensure the automated methods are transparent and positively build operator trust [6], [14].

The paper, thus, introduces a new AI-Augmented Inference Engine (RCASage) based on the principles of causal discovery, anomaly detection, and proactive defect discovery, to attain Autonomy (reducing human interactions), Explainability (with interpretable failure paths), and Scalability (that serves thousands of services). The paper presents a novel, vertical processing model that incorporates proactive software quality and reactive incident remediation, filling a significant gap in the current methods of anomaly detection in distributed systems. The technical novelty of the given approach is a particular combination of Graph-based Deep Learning (GBDL) and Probabilistic Causal Inference in a real-time processing feedback loop. As opposed to current methods that examine causal discovery as a "reactivity-focused" statistical procedure using a post-mortem or offline Granger Causality method such as in standard PC algorithms or offline Granger Causality analysis, RCASage instantiates a Live Service Dependency Graph, an ever-refreshing system that runs in synchronism with the Neural Granger Causal Discovery algorithm. This new system provides a much greater speed of assessment and segregation of propagated noise and initial faults in a sub-second latency range. RCASage can be used to form a new proactive closed-loop prognostic system in combination with Just-in-Time (JIT) methods to detect defects at the CI/CD level, allowing unprecedented levels of system troubleshooting, where code-level metadata guides system diagnostic-level troubleshooting. This study is clearly a step towards evolving the field beyond the limits of anomaly detection or reactive system paradigms (AIOps 1.0) into a novel intelligent system that seeks to explain identified anomalies and their interactions (AIOps 2.0), thus offering a solid, adaptable technical blueprint on how to solve intractable operational issues in contemporary DevOps and Site Reliability Engineering (SRE) practices [6], [10].

The evolution of RCA, the suggested architecture, and the mathematical representations of the RCASage inference engine, as well as the performance benchmarks that prove the approach to minimizing MTTR while providing human-readable causal chains, are explained in the following sections.

## 2        REVIEW OF RELATED WORKS

The historical development of Root Cause Analysis (RCA) in distributed systems evolved from manual and reactive processes to automated intelligence systems. This literature can be narrowed down to two (complementary) paradigms: Causal Inference-based Approaches and Graph-based Modelling, and a growing attention to the incorporation of proactive systems.

### 2.1        Evolution of RCA from Traditional Processes to AIOps

Traditionally, RCA was based on a manual analysis of logs and alerts, custom rule-based systems, or pre-defined dependency maps [6], [7], [15]. Generally, this approach has been pervasively deficient for contemporary

microservices topologies with dynamic systems, transient elements, and just the volume of heterogeneous telemetry data they generate [1], [2], [6], [7]. Figure 2.1 illustrates the traditional view of RCA. The diagram represents the traditional six-step lifecycle of RCA that is human-centric. This manual technique is being undermined by the intricacy and magnitude of contemporary distributed microservice systems, where fault detection and dependency mapping have become hard to handle effectively by human engineers [1], [2], [4].
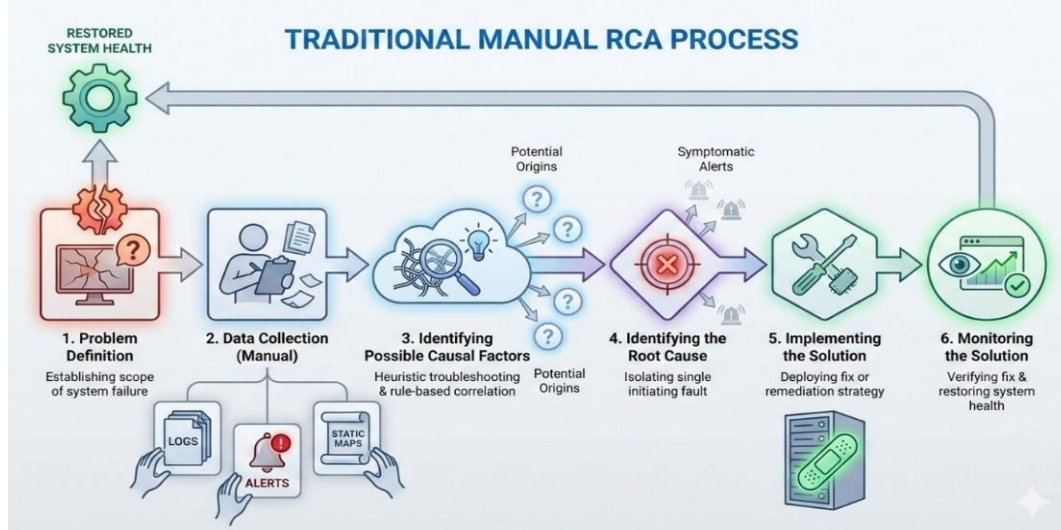


*Figure 2.1 Conceptual Framework of the Traditional Approaches to Root Cause Analysis. (Image Source: Researcher's Diagram drawn with the aid of Nano Banana Pro)*

Despite the significant step forward introduced by the concept of AI for IT Operations (AIOps), the utilization of machine learning algorithms to detect anomalies, analyse logs, and correlate events to mitigate alert fatigue, many of the tools are symptom-based, clustering the anomalies rather than directly inferring the cause and effect of the event [6], [7], [15]. Such a constraint reiterates a major necessity to move past correlation to actual causation, culminating in a body of work investigating alternative approaches to the issue, including causal AI and graph-based models, to capture the intricate dynamics of failure in complex systems [8], [11], [15].

2.2     Causal Inference for Root Cause Identification

The methods of causal inference directly address the "Observability-Complexity" gap by trying to decouple the causative events and the symptomatic correlations [6], [10], [16]. The inherent purpose here is to achieve a Causal AI-based Root Cause Identification that perceives failures as an intervention that empowers systems to uncover the causal graph underlying failures, using observational data [6], [7], [11]. Techniques, including Granger Causal Discovery, have been applied to time-series metrics to determine causal graphs in microservices [16], [17]. Probabilistic model techniques, including Bayesian networks, and probabilistic models such as Beta-Binomial, in industrial Root Cause Identification (RCI) engines target the concepts of differential observability to enable fault localization by finding

invariance in distributions under conditions of failure [6], [7]. They address the issue with more definitive solutions than correlation analysis, but they have serious problems with their scalability, flexibility, or even their ability to work in real-time without large volumes of supervised training data [7], [8].

2.3     Graph-Based and Multi-Modal Modelling

Graph-based models provide a natural representation paradigm to the inevitably interconnected and complex nature of microservice, with nodes (services, hosts, and events) and edges (dependencies, communications) [7], [9], [15]. This paradigm is very efficient in the dynamic relations of non-Euclidean data modelling. The research in this direction extends to metric-level RCA frameworks, which combine multi-modal system telemetry (logs, metrics, and distributed traces) to arrive at specific fault localization [18], up to highly advanced frameworks such as GROOT, which constructs real-time graphs of causality events such as fine-grained system monitoring events as nodes to make more precise inferences [6], [7]. Graph-Based Deep Learning (GBDL) has also been considered as an efficient method of anomaly detection as a valuable predecessor to RCA [9], [15]. These models are proficient at managing the propagation of failures among complicated dependency systems.

2.4     Integrating Proactive AI into the Software Lifecycle

A holistic AI-based RCA system is not only based on the concept of reactive diagnosis but also focuses on proactive defect and failure prediction throughout the software delivery lifecycle [3], [19]. This is achieved through the application of AI in the entire CI/CD pipeline to make independent decisions [20]. Specific applications of this are:

- Machine learning-based Automated Anomaly Detection in CI/CD Pipelines [21].
- Optimization of CI/CD Pipelines speed and reliability using AI [22], [23].
- Just-in-time (JIT) Software Defect Prediction (SDP), with deep learning, process/quality measures to predict defects, bugs, or commits that cause failures [24], [25], [26], [27], [28], [29], [30], [31], [32].
- Prediction of failure of cloud-native applications based on monitoring data [33].

A combination of these proactive features with the reactive causal-inference functions forms a whole and intelligent system with the ultimate aim of sustaining stability, and eventually transforming RCA from a human-centered, hypothesis-based process to an AI-based, evidence-driven science [6], [10], [11].

## 3    METHODOLOGY

The paper proposes the development of an AI-based augmented inference engine (which will be referred to as RCASage) that will have the ability to offer root cause analysis of microservices in a distributed system using both proactive defect prediction to detect errors in its execution and reactive causal analysis/reasoning. The approach utilizes a three-tier system shown in Figure 3.1, which consumes multi-modal telemetry, carries out anomaly scoring, and infers causation.
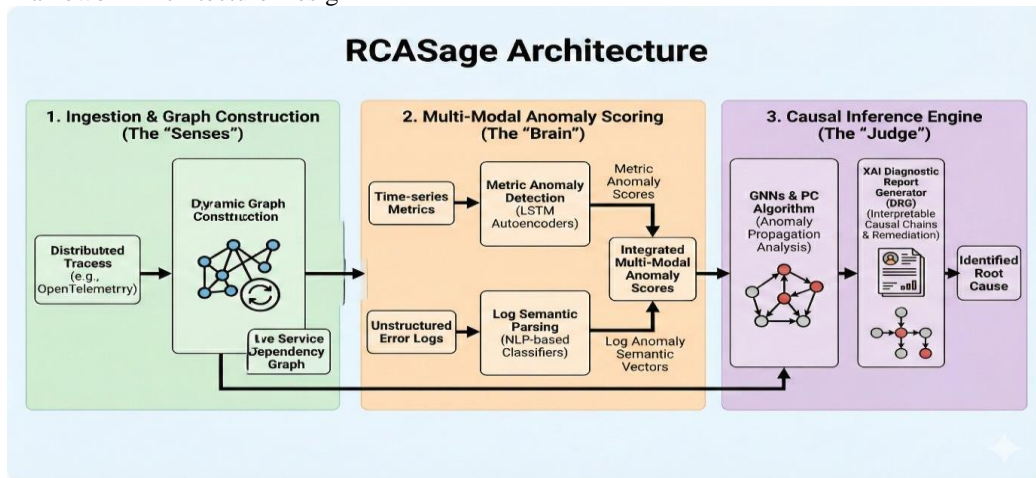
Framework Architecture Design



*Figure 3.1 System Architecture depicting the workflow and interactions of the proposed AI-Augmented Inference Engine – RCASage, sourced from the literature reviews of Chaudry [10], handling ingestion & graph building, Adenekan [14], handling multi-modal anomaly scoring, and Ikram et al. [16], and Lin et al. [17], handling the causal inference engine (Image Source: Researcher's Diagram drawn with the aid of Nano Banana Pro)*

### A.    Data Ingestion and Graph Construction (The Senses)

A Data Ingestion and Normalization Layer (DINL) is the first layer that would enable situational awareness. It has the task of combining all heterogeneous system telemetry, e.g., distributed traces (through OpenTelemetry), logs, time-series metrics (latency, CPU), and CI/CD metadata information (commit history, build output) [34], [35]. This information is then centralized and normalised towards a uniform model over time, ensuring accurate timestamps [7]. In addition to the mere aggregation of data, the DINL builds an Augmented Service Dependency Graph (ASDG), the nodes of which are related to the various components of the target system (servers, microservices). It is necessary to note that, unlike traditional

system maps, the ASDG represents dynamic functional dependencies among components as weighted edges that allow a structural "topological mask" that serves as the original search space of the causal engine itself, which guarantees that all the subsequent inferences are appropriately linked back to the physical reality of the target system [7], [15], [19]. The live graph is also continuously revised to reflect the current structure of the target system, which gives the required context of the structural aspect for all other analyses [6].

### B.    Multi-Modal Anomaly Detection and Scoring (The Brain)

The parallel and multi-modal anomaly detection processing layer is used to diagnose the system using raw telemetry data. An anomaly detector called the Real-Time Anomaly Detector (RTAD) uses both unsupervised learning algorithms and statistical models to detect any irregularities or outlier properties of the system behaviour [7], [9], [19]. It also involves using unsupervised LSTMs in the form of autoencoders to study time-series data, where crucial parameters are concerned [6] and using Isolation Forest algorithms to identify any anomaly/outlier in the logs [7]. At the same time, NLP classifiers analyse unstructured logs and extract semantic vectors in order to identify salient error patterns [6], [17]. The twist at this phase is the generation of a fused Multi-Modal Anomaly Vector ($V_a$) per service node. The framework also generates an Anomaly Significance Index (ASI) by summing up any variations in the parameters. This index has the benefit of not only alerting the user, but also as a

mathematical notion of pointing various nodes that are considered "suspicious" to the causal inference step, which basically removes any noise in the system before it is passed on to the discovery algorithms. The multi-stream method effectively handles the heterogeneity and the presence of an inherent class imbalance in the operational data to produce a homogenized set of timestamped anomaly scores to be used in causal analysis.

**C.        Causal Inference and Root Cause Identification (The Judge)**

The main innovation of the judge is the use of the Anomaly-Constrained Causal Discovery (AC-CD) algorithm in the reasoning process via the AIE. The engine uses a soft-intervention graph to define the system failures as soft-interventions, changing the state of a component [7]. Since the existing PC algorithms are categorically highly complex in high-dimensional microservice environments, RCASage presents a further step known as the Search-Space Pruning step to address the existing problems in the following ways:

1.        Topological Constraint: It makes use of the ASDG (Section A) to eliminate physically impossible causal edges.

2.        Heuristic Constraint: It uses the ASI (Section B) to initialize the PC algorithm with a skeleton using those nodes that have significantly strong anomalies.

3.        Refined Discovery: The engine then goes on to use Neural Granger Causal Discovery only on this restricted set.

To deal with uncertainty and real-world noisy data, it integrates Beta-Binomial inference (to explain beliefs in component health ($\theta_i \sim$ Beta ($\alpha_i$, $\beta_i$))) and represents updating those beliefs when a request is issued, with observed outcomes [7]. The engine combines Graph Neural Networks (GNNs) with the Neural Granger Causal Discovery algorithm for causal inference [6], [8], [15], [16], [17]. In this case, the GNNs consider propagation paths of anomalies on the graph, while the causal inference algorithms infer the structure based on the observation [18]. RCASage achieves computational and statistical robustness by limiting the searches of the PC algorithm through multi-modal scores, and can provide sub-second fault localization to uniquely distinguish between initiating errors and symptomatic behaviour due to cascading path consequences [6], [10], [11].

**D.        Proactive Integration and Explainable AI (XAI)**

To create a unified framework, proactive Software Defect Prediction (SDP) is included as a component of the CI/CD pipeline. The predictive models of JITLine and GHA-BFP use hybrid metrics and bi-modal learning to predict commits that cause defects or building failures, thus preventing defects before they reach production [28], [30], [31], [32]. The framework also uses a Causal Feedback Loop to trace the root cause found in production to a particular CI/CD metadata, thus continuously updating the SDP models. Finally, the framework includes a Diagnostic Report Generator (DRG) which emphasizes operational usability by introducing Explainable AI (XAI) methods to achieve interpretable results, in the shape of a ranked list of probable root causes with confidence scores attached, causal diagrams, and suggested recommendatory actions, turning "black-box" AI results into actionable, evidence-based insights, more likely to build trust and support rapid resolution [7], [14].

3.1        Key Mathematical Concepts of the Proposed Architecture

Essentially, algorithms or mathematical models employed in predicting, detecting anomalies, and causal inference are the strength of the proposed methodology, elevating the RCASage framework from correlation to causal inference.

**A.        Formal Graph Representation and GAT-based Anomaly Propagation:**

The system is modelled as a dynamic directed graph $G = (V, E, X)$, where $V$ represents the set of $N$ microservices [9], [15]. The Augmented Service Dependency Graph (ASDG) is formally derived from distributed traces such that an edge $e_{ij} \in E$ exists if a functional dependency or request flow is observed between service $s_i$ and $s_j$ [19], [21].

To model how anomalies propagate through this topology, RCASage utilizes Graph Attention Networks (GAT). Unlike standard GCNs, GAT allows for masked self-attention layers to assign different levels of importance to neighbouring services:

$$h_i' = \sigma\left(\sum_{j \in N_i} \alpha_{ij} W h_j\right)$$

Where:

•        $N_i$ is the neighbourhood of service $i$.

•        $\alpha_{ij}$ is the attention coefficient, representing the strength of the dependency/influence of the service $j$ on service $i$:

$$\alpha_{ij} = \frac{exp\left(LeakyReLU\left(a^T[Wh_i||Wh_j]\right)\right)}{\sum_{k \in N_i} exp\left(LeakyReLU\left(a^T[Wh_i||Wh_k]\right)\right)}$$

•        This formalization allows the engine to weight "noisy" dependencies lower than critical failure paths during the inference stage.

B.       Neural Granger Causal Discovery (NGCD):

To capture non-linear causal relationships in multi-modal telemetry time-series ($x_t$), RCASage implements Neural Granger Causal Discovery [8], [10], [16], [17]. It models the evolution of a metric $x^{(i)}$ for service $i$ as a non-linear function of the past values of all services $j \in \{1, ..., N\}$:

$$x_t^{(i)} \approx f_i\left(x_{<t}^{(1)}, x_{<t}^{(2)}, ..., x_{<t}^{N}\right)$$

In this framework, $f_i$ is a Component-wise Multilayer Perceptron (MLP). Service $j$ is determined to causally influence service $i$ if the weights of the first layer of $f_i$ corresponding to $x^{(j)}$ are non-zero. This is enforced during training via a group Lasso penalty to ensure a sparse causal graph:

$$\min_f \sum_t ||x_t - f(x_{<t})||^2 + \lambda \sum_{i=1}^{N} \sum_{j=1}^{N} ||W^{(ij)}||_2$$

Where $W^{(ij)}$ represents the weight matrix connecting service $j$'s history to service $i$'s prediction. If $||W^{(ij)}||_2 > 0$, a causal edge $j \rightarrow i$ is confirmed [18].

C.       Beta-Binomial Health Update Rule:

The Autonomous Inference Engine (AIE) quantifies the probability of a component being the root cause using Bayesian updates. The health $\theta_i$ of service $s_i$ is modelled as a Beta distribution $Beta(\alpha_i, \beta_i)$ [7]. Upon observing $n$ requests with $k$ successes, the posterior health is updated via the conjugate prior rule:

$$P(\theta_i|data) = Beta(\alpha_{prior} + k, \beta_{prior} + (n - k))$$

The Root Cause Score ($S_{rc}$) for a node is then calculated as the negative log-likelihood of the observed failure rate given the current health belief, prioritizing nodes whose performance deviates most sharply from their historical Bayesian baseline.

D.       Just-In-Time (JIT) Defect Prediction Model:

For proactive prevention, the framework employs a deep-learning-based classifier for JIT Defect Prediction [29], [36]. Given a code change $c$ the model predicts the probability of a defect $P(Y_c = 1)$ by fusing a feature vector $\phi(c)$ consisting of Delta Metrics (e.g., additions/deletions of code lines, author expertise) and Semantic Change Embeddings:

$$P(Y_c = 1|\phi(c)) = Sigmoid(W_2 \cdot ReLU(W_1\phi(c) + b_1) + b_2)$$

This formalization can identify "bug-inducing commits" with great precision and assists the system in indicating possible production-grade problems prior to the deployment of the code [25], [26], [27], [30], [31], [32].

## 4        PERFORMANCE BENCHMARKS AND DISCUSSIONS

The proposed AI-Augmented Inference Engine (AAIE) model can offer material changes in the speed, precision, and practice in the diagnostic process of microservices ecosystems. This discussion is a synthesis of the anticipated results of the architectural analysis and empirical standards that are established in the corresponding literature.

4.1        Quantitative Gains in Diagnostic Accuracy and Speed

Empirical research and the implementation of similar AIOps systems demonstrate that AI-based RCA can be much more effective than classical algorithms. The Accuracy of the Root Cause Identification is a vital KPI; thus, its level is likely to increase significantly. For example, hybrid graphical models have reached accuracy rates of more than 90% already, which is already a substantial improvement over the more traditional system monitoring techniques (~57%), and the average AIOps tools (~74%) [6], [7], [18]. The "root cause" vs. "propagated effects" dilemma of distributed environments can be directly tackled due to such high accuracy rates, establishing that the framework transcends correlation-based analysis to engage in causal analysis [8], [10], [11]. The time-to-resolution is the most significant operational impact. The proposed framework eliminates the manual traceback processes with the automation of multi-modal telemetry and dependency graphs analyses. According to the results of empirical studies conducted on similar work, Mean Time to Resolve (MTTR) and Time to Diagnosis (TTD) can be reduced by more than 90 percent in comparison to conventional approaches [6], [13]. End-to-end diagnosis times of less than 5 seconds are achievable based on production-level studies on similar systems. Inference latencies are also kept low (e.g. below 900ms) even when the event loads are high [6], [7]. The typical figures and percentages of the performance gains (in time) of AI-Augmented and Causal AI systems over the traditional, manual, or correlation-based approaches are shown in Table 4.1.

*Table 4.1 Comparative Performance Gains of Root Cause Analysis Methods*

| Metric | Traditional RCA (Manual/Correlation-based) | AI-Driven Framework (Causal/Graph-based) | Improvement (%) |
|---|---|---|---|
| Mean Time To Resolution (MTTR) | Hours to Days | Minutes to Hours | 30% to 78% Reduction [4], [12], [13], [37] |

| Mean Time To Detect (MTTD) | Hours | Seconds to Minutes | 42% Improvement in RCA Identification Time [4] |
|---|---|---|---|
| RCA Accuracy (Correct Root Cause Identified) | ~ 60\% - 78\% | ~ 90\% - 95\% | 17% to 60% Boost in Accuracy [7], [20] |
| False Positives/Alert Noise | High | Low | 50% to 80% Fewer False Positives [37] |

This acceleration is due to the inclusion of the Graph-Based Deep Learning module for holistic anomaly detection and the Neural Granger Causal Discovery module for causal analysis [9], [17], [19].

### 4.2 Addressing Operational Complexity and Proactive Stability

The framework is aimed at addressing the complexity of the high number of services contained within a distributed microservices setting. The dynamic Service Dependency Graph allows the necessary architectural context for the analysis of the failures distributed across thousands of services [15], [19]. Moreover, the multi-modal data (metrics, logs, and traces) fusion is required for precise fault investigation and detecting faults in a heterogeneous setup, getting around the constraints of single-data-type analysis [18].

A key milestone that has been implemented is the shift towards a proactive mode, as opposed to reactive modes of operation. The support of Just-in-Time (JIT) defect prediction in the CI/CD pipeline allows early detection and correction of commits that cause bugs in a manner that prevents incidents in production in the first place [27], [30], [31], [37]. Such proactive measure along with the ability to anticipate the failure of cloud-native applications based on the monitoring data is a crucial preventive barrier [3], [20], [21], [33].

### 4.3 Operationalizing Causal AI: From Theory to Practice

The major advantage of the RCASage framework is that Causal AI is operationalized in high-velocity production conditions. Even though the underlying theoretical advantages of the method have been well-established, it has proven challenging to scale in real-time conditions [8], [11]. The proposed system avoids this by using an Anomaly-Constrained Causal Discovery (AC-CD) pipeline that makes sure that the significant computation of the causal engine is only applied to the topographically meaningful and statistically anomalous subsamples of services.

#### 4.3.1 The Role of the Diagnostic Report Generator (DRG)

The most important "last mile" component is the Diagnostic Report Generator (DRG), which presents an organized and readable format for abstract mathematical probability calculations of algorithms. The DRG component also acts as a Virtual Tier-1 Analyst, unlike the normal dashboarding tools, which only act as a viewer of telemetry information, the DRG component presents a multi-layered diagnostic report for:

1. Visual Causal Traversal: This will enable it to visualize the ASDG (Augmented Service Dependency Graph) and highlight the "hot paths," demonstrating how a database-level failure could propagate through the middleware to the end-user API.

2. Evidence Attribution: Each RCI has some marker indicators of evidence ("evidence tokens"), e.g., a specific pattern of log errors or a dramatic change in the Neural Granger weight, to tell why this particular service was the one being flagged among the neighbouring services.

3. Remediation Mapping: It associates the detected issue with the history of remediation and CI/CD metadata to offer suggested measures (e.g., "Roll back commit $x$" or "Increase thread pool for service $y$").

#### 4.3.2 Incorporating Trust in SRE with the embedding of Explainable AI (XAI)

The principles of Causal-centric XAI are incorporated into the proposed framework to make the prediction approach using the black-box model more acceptable during the operation phase. This is achieved in the proposed structure through the assistance of the following two specific mechanisms:

- Counterfactual Reasoning: The XAI layer enables SREs the capability of conducting a "What-if" analysis via the UI (e.g., "Would this latency exist in a case where Service A was healthy?). These soft interventions are simulated so that the framework produces contrastive explanations that are assuring to the engineers that the root cause diagnosis is correct.

- Path-Based Transparency: The framework also provides transparency using the concept of "justification paths" rather than giving a generic probability score. It also visualizes the Search-Space Pruning process to give an understanding of what services were not considered in the evaluation and what was the rationale behind such decisions, therefore, breaking the barrier of scepticism when deploying AI in automation.

Such degrees of transparency are necessary to win the confidence of SRE teams. Table 4.2 [21], [26] shows the performance of different ML/DL models integrated in the CI/CD pipeline to predict defects and failures. The F1-score and Area Under the Curve (AUC) are actually typical in defect prediction systems to perform classification tasks. Nevertheless, as it has been observed in some recent literature [6], [7], [14] high accuracy without explanations leads to "alert skepticism." The operationalization of XAI within RCASage guarantees that the

framework can be used as an operational companion to human teams so that they can establish a feedback loop of continuous learning to refine future causal issues.

*Table 4.2 Table showing the performance of Proactive Prediction Models*

| Prediction Task | Model Used (e.g., JITLine, GHA-BFP, Random Forest) [28], [30], [36] | Key Metrics Used [32] | Performance Metric (e.g., F1-Score) | Benefit (e.g., Reduced Build Failures) [28], [37] |
|---|---|---|---|---|
| JIT Defect Prediction [30] [31] | Bi-modal Change Representation Learning | Change Metrics, Source Code Metrics | $F1_{JIT}$ | Code Quality, Reduced Production Incidents |
| Automated Anomaly Detection in CI/CD [21] | Machine Learning Approaches | Build Logs, Test Results | $AUC_{Anomaly}$ | Early Fault Isolation, Faster Feedback |
| Build Failure Prediction [28] | GHA-BFP | Repository Metrics, Build Performance Metrics [35] | Accuracy | CI/CD Pipeline Optimization [23] |

## 5    CONCLUSION AND FUTURE WORK

### 5.1    Conclusion

This study proposes an architectural design for RCASage, an AI-augmented inference engine (AAIE) conceived as a holistic conceptual framework for addressing critical challenges in root cause analysis (RCA) within distributed microservices environments. The proposed framework advances beyond reactive troubleshooting by introducing an automated, evidence-based approach through a vertically integrated pipeline that combines dynamic dependency modeling, multi-modal anomaly detection, and anomaly-constrained causal discovery.

The design is specifically intended to enable a transition from symptomatic correlation to explicit cause-and-effect reasoning, supported by Explainable AI (XAI) mechanisms that ensure transparency and foster operational trust. The framework adopts a dual focus by integrating proactive defect prediction within the CI/CD pipeline alongside a high-precision reactive RCA engine, thereby offering a comprehensive approach to system reliability.

Although full empirical evaluation of the RCASage architecture remains the subject of forthcoming experimental research, existing literature on comparable causal-graph-based architectures and autonomous inference agents provides insight into the performance potential of such systems. Drawing on evidence reported in related studies, estimated benchmarks suggest the potential to reduce Mean Time to Resolution (MTTR) and Time to Diagnosis (TTD) by more than 90% in fully autonomous environments. This work therefore represents a foundational step toward the development of highly resilient, self-adaptive, and ultimately autonomous IT operations.

### 5.2    Future Work and Research Directions

Despite the theoretical advancements of the framework, various research areas need to be explored to develop this blueprint into a production-ready system:

- Experimental Validation and Prototyping: The next step to take would be to implement the RCASage pipeline into a physical controlled microservices environment (like Google's Online Boutique or other similar benchmarking suites) to provide primary data and test the performance of the AC-CD algorithm.

- Improved Learning and Adaptation: With the addition of a Reinforcement Learning layer, the system would be able to learn from the results of remediation, hence improving its inference methods over time. In addition, it would be beneficial to use direct feedback provided by Site Reliability Engineers to raise the threshold of anomalies and dependency scores over time.

- Advanced Modelling and Generalization: Agent-Based Modelling studies about digital twin infrastructure have the potential to improve the prediction methods by modelling the dynamics of error propagation. Also, work on the generalization abilities of basic architectural ideas to other complex domains, including healthcare implementation or Industrial IoT, will be an interesting path forward.

- Technical Refinements and Benchmarks: Future work should involve work to overcome such limitations as the "cold start" problem, and the use of more sophisticated forms of telemetry (e.g., execution metrics, or logs of developer activity). Standard benchmarks for evaluating the scalability and performance of Causal AI engines for ultra-large-scale deployments would also form another significant area that contributes to the advancement of the field.

- Fully Autonomous Architectures: The future direction will be towards fully autonomous systems. This requires studying how to remove the requirement of manually defined failure times, advancing towards end-to-end anomaly detection and RCA that enable self-healing cloud architectures.

## REFERENCES

[1]     A. Das, "5 Common Challenges in Large-Scale Microservices," Medium, 12 March 2025. [Online]. Available: https://article.arunangshudas.com/5-common-challenges-in-large-scale-microservices-c199332ad571. [Accessed 24 November 2025].

[2]     M. Söylemez, B. Tekinerdogan and A. Kolukısa, "Challenges and Solution Directions of Microservice Architectures: A Systematic Literature Review," *Applied Sciences,* vol. 12, no. 11, p. 5507, 29 May 2022.

[3]     A. Sappa, "Adaptive Machine Learning Models for Effort Estimation and Risk Detection in Distributed Software Project Pipelines," *Journal of Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications (JoWUA),* vol. 16, no. 2, pp. 54-74, June 2025.

[4]     Cognizant 20-20 Insights, "Pitfalls & Challenges Faced During a Microservices Architecture Implementation," Cognizant Digital Systems & Technology, 2020.

[5]     Z. Liu, G. Fan, H. Yu and L. Chen, "An Approach to Modeling and Analyzing Reliability for Microservice-Oriented Cloud Applications," *Wireless Communications and Mobile Computing,* vol. 2021, no. 1, pp. 1-17, August 2021.

[6]     H. Wang, P. Nguyen, J. Li, S. Kopru, G. Zhang, S. Katariya and S. Ben-Romdhane, "GRANO: interactive graph-based root cause analysis for cloud-native distributed data platform," *Proceedings of the VLDB Endowment,* vol. 12, no. 12, pp. 1942-1945, 1 August 2019.

[7]     H. Wang, Z. Wu, H. Jiang, Y. Huang, J. Wang, S. Kopru and T. Xie, "Groot: An Event-graph-based Approach for Root Cause Analysis in Industrial Settings," in *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, Melbourne, Australia, 2021.

[8]     L. Pham, H. Ha and H. Zhang, "Root Cause Analysis for Microservice System based on Causal Inference: How Far Are We?," in *ASE '24: 39th IEEE/ACM International Conference on Automated Software Engineering*, Sacramento, CA, USA, 2024.

[9]     A. D. Pazho, G. A. Noghre, A. A. Purkayastha, J. Vempati, O. Martin and H. Tabkhi, "A Survey of Graph-Based Deep Learning for Anomaly Detection in Distributed Systems," *IEEE Transactions on Knowledge and Data Engineering,* vol. 36, no. 1, pp. 1-20, 2024.

[10]    A. Chaudhry, "AI-Driven Root Cause Analysis: Leveraging Causal Inference in SRE," Medium, 7 January 2025. [Online]. Available: https://amitchaudhry.medium.com/ai-driven-root-cause-analysis-leveraging-causal-inference-in-sre-228c236e4313.

[11]    S. Jha, A. Rahane, L. Shwartz, M. Palaci-Olgun, F. Bagehorn, J. Rios, D. Stingaciu, R. Kattinakere and D. Banerjee, "Causal AI-based Root Cause Identification: Research to Practice at Scale," p. 2025.

[12]    O. Israel, "AI-Augmented Root Cause Analysis for System Failures in Cloud and Edge Data Centers," September 2025.

[13]    B. T. Tutuncuoglu, "AI-Augmented Root Cause Analysis: Autonomous Inference Engine for Multi-Layered System Outages," *IEEE,* pp. 1-18, 2025.

[14]    T. K. Adenekan, "Explainable AI Techniques for Root Cause Analysis in Complex Systems," pp. 1-14, 2024.

[15]    Á. Brandón, M. Solé, A. Huélamo, D. Solans, M. S. Pérez and V. Muntés-Mulero, "Graph-based root cause analysis for service-oriented and microservice architectures," *Journal of Systems and Software,* vol. 159, p. 110432, 2020.

[16]    A. Ikram, S. Chakraborty, S. Mitra, S. Saini, S. Bagchi and M. Kocaoglu, "Root Cause Analysis of Failures in Microservices through Causal Discovery," *Advances in Neural Information Processing Systems,* vol. 35, pp. 31158-31170, 2022.

[17]    C.-M. Lin, C. Chang, W.-Y. Wang, K.-D. Wang and W.-C. Peng, "Root Cause Analysis in Microservice Using Neural Granger Causal Discovery," *Proceedings of the AAAI Conference on Artificial Intelligence,* vol. 38, no. 1, pp. 206-213, 25 March 2024.

[18]    Y. Wang, Z. Zhu, Q. Fu, Y. Ma and P. He, "MRCA: Metric-level Root Cause Analysis for Microservices via Multi-Modal Data," in *ASE '24: Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, 2024.

[19]    I. Erakovic and C. Pahl, "Root Cause Analysis for Microservices based on Architecture Profiling and Anomaly Detection," in *Cloud Computing and Services Science*, Porto, Portugal, Springer-Verlag, 2025, pp. 1-29.

[20]    M. Baqar, S. Naqvi and R. Khanda, "AI-Augmented CI/CD Pipelines: From Code Commit to Production with Autonomous Decisions," pp. 1-13, 16 August 2025.

[21]    F. Williams, A. A. Adeniran, E. Edmund and B. Idowu, "Automated Anomaly Detection in CI/CD Pipelines Using Machine Learning Approaches," *ResearchGate,* pp. 1-15, 1 September 2025.

[22]    R. C. Thota, "CI/CD Pipeline Optimization: Enhancing Deployment Speed and Reliability with AI and Github Actions," *International Journal of Innovative Research in Engineering & Multidisciplinary Physical Sciences (IJIRMPS),* vol. 8, no. 2, pp. 1-11, March 2020.

[23]    R. Farihane, I. Chlioui and M. Radgui, "CI/CD Pipeline Optimization Using AI: A Systematic Mapping Study," *Engineering Proceedings,* vol. 112, no. 1, 9-11 July 2025.

[24]    E. N. Akimova, A. Y. Bersenev, A. A. Deikov, K. S. Kobylkin, A. V. Konygin, I. P. Mezentsev and V. E. Misilov, "A Survey on Software Defect Prediction Using Deep Learning," *Mathematics,* vol. 8, no. 1180, pp. 1-14, 2021.

[25]    Q. Yu, S. Jiang, J. Qian, L. Bo, L. Jiang and G. Zhang, "Process metrics for software defect prediction in object-oriented programs," *IET Software,* vol. 14, no. 3, pp. 283-292, 1 June 2020.

[26]    Durga and D. A. Sinha, "Enhancing Software Reliability through Intelligent Fault Prediction Using Machine Learning," *International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT),* vol. 11, no. 3, pp. 945-956, 11 June 2025.

[27]    A. B. Tavakkoli, "Framework for Bug Inducing Commit Prediction Using Quality Metrics," Electronic Thesis and Dissertation Repository, Ontario, Canada, 2024.

[28]    J. Li, Y. Zhang, T. Wang and Y. Wu, "GHA-BFP: Framework for Automated Build Failure Prediction in GitHub Actions," in *2024 31st Asia-Pacific Software Engineering Conference (APSEC)*, Chongqing, China, 2024.

[29]    M. J. Haruna and T. C. Darius, "Software Defect Prediction Using Machine Learning and Deep Learning Techniques," *Kasu Journal of Computer Science,* vol. 1, no. 3, pp. 527-543, September 2024.

[30]    C. Pornprasit and C. K. Tantithamthavorn, "JITLine: A Simpler, Better, Faster, Finer-grained Just-In-Time Defect Prediction," in *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*, Madrid, Spain, 2021.

[31]    Y. Jiang, B. Shen and X. Gu, "Just-In-Time Software Defect Prediction via Bi-modal Change Representation Learning," *The Journal of Systems & Software,* vol. 219, no. C, pp. 1-14, 1 January 2025.

[32]    A. Zaim, J. Ahmad, N. H. Zakaria, G. E. Su and H. Amnur, "Software Defect Prediction Framework Using Hybrid Software Metric," *International Journal on Informatics Visualization,* vol. 6, no. 4, pp. 921-930, December 2022.

[33]    L. Toka, G. Dobreff, D. Haja and M. Szalay, "Predicting cloud-native application failures based on monitoring data of cloud infrastructure," in *2021 IFIP/IEEE International Symposium on Integrated Network Management (IM)*, Bordeaux, France, 2021.

[34]    A. Houerbi, C. Siala, A. Tucker, D. E. Rzig and F. Hassan, "Empirical Analysis on CI/CD Pipeline Evolution in Machine Learning Projects," in *Association for Computing Machinery (ACM) Conference (Conference'17)*, New York, NY, USA, 2025.

[35]    R. S. Constantinescu, "Exploring Descriptive Metrics of Build Performance: A Study of GitHub Actions in Continuous Integration Projects," 2023.

[36]    W. H. Seow, C. Y. Lim and S. L. Ang, "Random Forest Model for Software Build Time Prediction on CI/CD Pipeline," *Pertanika Journal of Science & Technology,* vol. 33, no. 2, pp. 1031-1048, 21 February 2025.

[37]    T. Rausch, W. Hummer, P. Leitnery and S. Schulte, "An Empirical Analysis of Build Failures in the Continuous Integration Workflows of Java-Based Open-Source Software," in *14th International Conference on Mining Software Repositories (MSR 2017)*, Buenos Aires, Argentina, 2017.