

MicroHECL: High-Efficient Root Cause Localization in Large-Scale Microservice Systems

Dewei Liu
Fudan University
Shanghai, China

Chuan He
Fudan University
Shanghai, China

Xin Peng
Fudan University
Shanghai, China

Fan Lin
Alibaba Group
Hangzhou, China

Chenxi Zhang
Fudan University
Shanghai, China

Shengfang Gong
Alibaba Group
Hangzhou, China

Ziang Li
Alibaba Group
Hangzhou, China

Jiayu Ou
Alibaba Group
Hangzhou, China

Zheshun Wu
Alibaba Group
Hangzhou, China

Abstract—Availability issues of industrial microservice systems (e.g., drop of successfully placed orders and processed transactions) directly affect the running of the business. These issues are usually caused by various types of service anomalies which propagate along service dependencies. Accurate and high-efficient root cause localization is thus a critical challenge for large-scale industrial microservice systems. Existing approaches use service dependency graph based analysis techniques to automatically locate root causes. However, these approaches are limited due to their inaccurate detection of service anomalies and inefficient traversing of service dependency graph. In this paper, we propose a high-efficient root cause localization approach for availability issues of microservice systems, called MicroHECL. Based on a dynamically constructed service call graph, MicroHECL analyzes possible anomaly propagation chains, and ranks candidate root causes based on correlation analysis. We combine machine learning and statistical methods and design customized models for the detection of different types of service anomalies (i.e., performance, reliability, traffic). To improve the efficiency, we adopt a pruning strategy to eliminate irrelevant service calls in anomaly propagation chain analysis. Experimental studies show that MicroHECL significantly outperforms two state-of-the-art baseline approaches in terms of both accuracy and efficiency. MicroHECL has been used in Alibaba and achieves a top-3 hit ratio of 68% with root cause localization time reduced from 30 minutes to 5 minutes.

Index Terms—microservice, availability, root cause localization, anomaly detection, service call graph

I. INTRODUCTION

Microservice architecture has been the latest trend in building cloud-native applications and more and more companies have chosen to migrate from the so-called monolithic architecture to microservice architecture [1]–[3]. Industrial microservice systems can include hundreds to thousands of services. For example, the e-commerce system of Alibaba contains more than 30,000 services and includes over 300 subsystems. Microservice systems support the core business of the companies, especially for Internet companies which provide online services through the Internet.

Availability issues have been a critical challenge for large-scale industrial microservice systems. These systems are highly dynamic and complex. A service can have several to thousands of instances running on different containers and

dynamically created or destroyed according to the scaling requirements at runtime; the execution of a microservice system may involve a huge number of microservice interactions and most of them are asynchronous and involve complex invocation chains [2]. In these systems, any anomaly with the quality (e.g., performance, reliability) of a service may propagate along service call chains and finally cause availability issues at the business level (e.g., drop of successfully placed orders). When an availability issue is detected, its root cause service and anomaly type need to be located in a short time (e.g., 3 minutes) to allow the developers to fix the issue quickly.

Existing approaches cannot efficiently support the root cause localization of availability issues for large-scale microservice systems. Developers in industry often rely on visualization tools to analyze logs and traces to identify possible anomalies and anomaly propagation chains to locate root causes [2], [3]. Researchers have proposed approaches for automatic root cause localization for microservice or the broader service-based systems using trace analysis [4], [5] or service dependency graph [6]–[12]. Trace analysis based approaches require expensive collection and processing of trace data, thus cannot efficiently work for large-scale systems. Service dependency graph based approaches construct service dependency graphs based on service calls and causal relationships (e.g., services co-located in the same machines). These approaches locate root causes by traversing service dependency graphs and detecting possible anomalies with the quality metrics (e.g., response time) of services. These approaches are limited due to their inaccurate detection of service anomalies and inefficient traversing of service dependency graphs, especially when the system has many services and dependencies.

In this paper, we propose a high-efficient root cause localization approach for availability issues of microservice systems, called MicroHECL. Given an availability issue initially reported on a service (called *initial anomalous service*), MicroHECL locates the root cause services and anomaly types (e.g., performance anomaly, reliability anomaly, and traffic anomaly) that cause the issue. It dynamically constructs a service call graph of the target microservice system, which reflects the service dependencies and related quality metrics

(e.g., response time) in the latest time window. Based on the graph, it analyzes possible anomaly propagation chains from the initial anomalous service by traversing the graph along anomalous service calls. To detect possible anomalies with individual service calls, we design an anomaly detection model for each of the three popular anomaly types (i.e., performance, reliability, traffic). To improve the efficiency of anomaly propagation chain analysis, MicroHECL adopts a pruning strategy to eliminate anomalous service call edges that are irrelevant to the current anomaly propagation chain. Finally, MicroHECL ranks candidate root causes based on the correlations between their quality metrics and the business metrics of the initial anomalous service.

To evaluate the effectiveness and efficiency of MicroHECL, we conduct a series of experimental studies with the data and availability issues collected from the e-commerce system of Alibaba. The results show that MicroHECL significantly outperforms two state-of-the-art baseline approaches Monitor-Rank [6] and Microscope [12] in terms of both the accuracy and efficiency. The results also confirm the effectiveness of the pruning strategy, which can significantly improve the efficiency while keeping the accuracy.

MicroHECL has been deployed in Alibaba for more than 5 months and used to handle more than 600 availability issues. It achieves a top-3 hit ratio of 68% and reduces the time of typical root cause localization from 30 minutes to 5 minutes. Feedback from the developers shows that most of them trust the recommendations of MicroHECL and treat the recommendations as the root causes by default. The analysis of the results also suggests some improvements of the approach.

The remainder of this paper is organized as follows. Section II introduces background knowledge about the problem scenarios. Section III presents an overview of MicroHECL and Section IV details the anomaly propagation chain analysis. Section V describes our experimental study for the evaluation of MicroHECL. Section VI introduces the practical application of MicroHECL in Alibaba. Section VII discusses related work and Section VIII concludes the paper.

II. BACKGROUND

The e-commerce system of Alibaba has more than 846 millions monthly active users. It adopts the microservice architecture and contains more than 30,000 services. These microservices are deployed with containers (i.e., Docker [13]) and virtual machines and orchestrated by the customized orchestrator based on Kubernetes [14], with Service Mesh [15] applied for service communications. The system is equipped with a large-scale monitoring infrastructure called EagleEye [16]. EagleEye includes a tracing SDK, a real-time ETL (Extract-Transform-Load) data processing system, a batch processing computing cluster, and a web-based user interface.

As the carrier of the business, the system needs to ensure high availability. Therefore, a business monitoring platform is deployed to raise timely alarms about availability issues. These availability issues usually indicate problems with the running of business, for example the drop of successfully placed orders

and success rate of transactions. An availability issue can be caused by different types of anomalies, each of which is indicated by a set of metrics. An anomaly can originate from a service and propagate along service calls, and finally cause an availability issue. In this work, we focus on the following three types of anomalies that cause most of the availability issues in Alibaba.

- **Performance Anomaly.** Performance anomaly is indicated by anomalous increase of response time (RT). It is usually caused by problematic implementation or improper environmental configurations (e.g., CPU/memory configurations of containers and virtual machines).
- **Reliability Anomaly.** Reliability anomaly is indicated by anomalous increase of error counts (EC), i.e., the numbers of service call failures. It is usually caused by exceptions due to code defects or environmental failures (e.g., server or network outage).
- **Traffic Anomaly.** Traffic anomaly is indicated by anomalous increase or decrease of queries per second (QPS). Anomalous traffic increase may break the services, while anomalous decrease may indicate that many requests cannot reach the services. Traffic anomaly is usually caused by improper traffic configurations (e.g., the traffic limits of Nginx [17]), DoS attack, or unanticipated stress test.

In anomaly detection, the 3-sigma rule is often used to identify outliers of metric values as candidates of anomalies. The 3-sigma rule states that in a normal distribution almost all the values remain within three standard deviations of the mean. The range can be represented as $(\mu - 3\sigma, \mu + 3\sigma)$, where μ is the mean and σ is the standard deviation. The values within the range account for 99.73% of all the values and the others can be regarded as outliers.

III. MICROHECL OVERVIEW

MicroHECL is a high-efficient root cause localization approach for availability issues. An availability issue is usually detected from the business perspective, e.g., the decline of successful orders. It may be caused by different types of anomalies. Currently MicroHECL supports three types of anomalies, i.e., performance anomaly, reliability anomaly, and traffic anomaly. When handling a specific type of anomalies, MicroHECL considers a corresponding service call metric and a specific direction of anomaly propagation. For example, when handling performance anomaly anomalies, MicroHECL considers response time of service calls and the anomaly propagation direction from downstream to upstream.

An overview of MicroHECL is presented in Figure 1. The approach includes the following three parts.

1) Service Call Graph Construction

When the runtime monitor detects an availability issue, MicroHECL starts a root cause analysis process. It constructs a service call graph of the target microservice system based on the service calls and metrics captured by the runtime monitor. The service call graph provides an updated snapshot of the microservice system by recording the service calls occurring

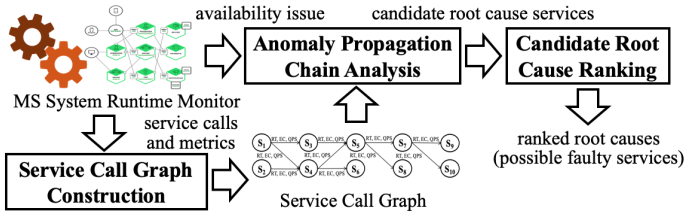


Fig. 1. MicroHECL Overview

in a latest time window (e.g., the last 30 minutes). To facilitate the root cause analysis of availability issues, the service call graph also records various quality metrics of service calls, such as response time (RT), error count (EC), and queries per second (QPS).

A service call graph consists a set of service nodes $S = \{S_1, S_2, \dots, S_n\}$. A node represents a service that is called in the latest time window (e.g., the last 30 minutes). Note that a service may have many instances in the running system. An edge $S_i \rightarrow S_j$ represents a service call from S_i to S_j that occurs in the latest time window. For each service call, the graph records a series of values for different quality metrics (e.g., RT, EC, QPS) in the latest time window (e.g., the last 60 minutes). For example, for each service call the graph records 60 RT values, each indicating the average response time of a minute.

The monitoring data (including service calls and their metrics) is stored in a time series database (TSDB). When an availability issue is raised, MicroHECL dynamically constructs a service call graph based on the data. Services and service calls in the graph are aggregated from service instances and calls between them. The metrics of a service call $S_i \rightarrow S_j$ are also aggregated from the corresponding metrics of the calls from the instances of S_i to those of S_j . It uses Flink [18], a distributed streaming data-flow engine, to aggregate the data. To optimize the performance, the service call graph is constructed in an on-demand and incremental way with the anomaly propagation chain analysis process: only when a service call is reached in the analysis the data about it is pulled from the TSDB.

2) Anomaly Propagation Chain Analysis

An availability issue is initially reported on a service (called *initial anomalous service*), but the root cause often lies in some of its downstream or upstream services. A root cause service and an initial anomalous service are usually connected by an anomaly propagation chain composed of a series of anomalous services. The purpose of anomaly propagation chain analysis is to identify a set of candidate root cause services by analyzing possible anomaly propagation chains from the initial anomalous service.

Based on the service call graph, MicroHECL analyzes possible anomaly propagation chains from the initial anomalous service of the availability issue. The analysis is done by traversing the service call graph along anomalous service call edges, starting from the initial anomalous service and

following the opposite directions of possible anomaly propagation directions. Each anomaly propagation chain considers a specific anomaly type propagated from a neighboring service of the initial anomalous service. To improve the efficiency, MicroHECL adopts a pruning strategy to eliminate anomalous service call edges that are irrelevant to the current anomaly propagation chain. This anomaly propagation chain analysis ends with a set of services as candidate root causes, each associated with a specific anomaly type. The analysis process together with the service anomaly detection method and the pruning strategy are detailed in Section IV.

3) Candidate Root Cause Ranking

MicroHECL ranks candidate root causes based on the possibility of causing the given availability issue. Based on our analysis of the monitoring data, we find that the anomaly index of the initial anomalous service has similar change trends with the anomaly index of the root cause service. Note that the anomaly index of the initial anomalous service is some kind of business metric (e.g., number of successful orders), while the anomaly index of a candidate root cause is some kind of quality metric (e.g., RT, EC, QPS). Therefore, we use the Pearson correlation coefficient [19] to measure the similarity of the change trends of the two anomaly indexes and rank the candidate root causes by the correlation coefficient.

For the initial anomalous service, the service call graph records a business metric value (e.g., number of successful orders) per minute. For a candidate root cause service, the service call graph records a value for the quality metric (e.g., RT, EC, QPS) of the corresponding anomaly type per minute. To reflect the latest change trends, we only consider the metric values in a latest time window (e.g., the last 60 minutes). The values for a specific business metric of the initial anomalous service in the last n minutes thus can be represented by a vector X , where X_i ($1 \leq i \leq n$) represents the value of the i th minute. Similarly, the values for a specific quality metric of a candidate root cause service in the last n minutes can be represented by a vector Y . Their Pearson correlation coefficient can be calculated as Equation 1, where \bar{X} and \bar{Y} represent the average of X and Y respectively. A positive (negative) value of $P(X, Y)$ indicates a positive (negative) correlation; and the larger the absolute value of $P(X, Y)$ (i.e., $|P(X, Y)|$), the more likely that the two change trends correlate. Therefore, we rank the candidate root causes by the absolute value of the Pearson correlation coefficient.

$$P(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}} \quad (1)$$

IV. ANOMALY PROPAGATION CHAIN ANALYSIS

In this section, we first introduce the process of anomaly propagation chain analysis, and then detail the pruning strategy and service anomaly detection methods.

A. Analysis Process

An availability issue may be caused by different types of anomalies. For different anomaly types the quality metrics

TABLE I
METRICS AND PROPAGATION DIRECTION OF AVAILABILITY ISSUES

Anomaly Type	Metric	Propagation Direction
Performance Anomaly	RT	downstream \rightarrow upstream
Reliability Anomaly	EC	downstream \rightarrow upstream
Traffic Anomaly	QPS	upstream \rightarrow downstream

that are considered in service anomaly detection and the anomaly propagation directions that are considered in propagation chain analysis are different. As shown in Table I, the considered quality metrics for performance anomaly, reliability anomaly, and traffic anomaly are response time (RT), error count (EC), and queries per second (QPS), respectively. The propagation directions for both performance anomaly and reliability anomaly are from downstream to upstream, while the propagation direction for traffic anomaly is from upstream to downstream.

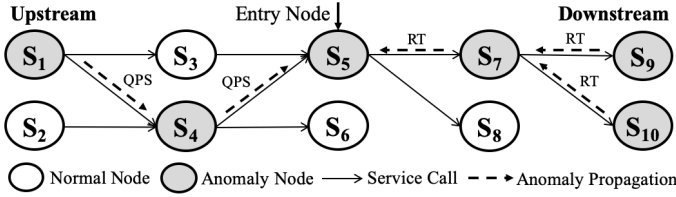


Fig. 2. Anomaly Propagation Chain Analysis Process

The process of the anomaly propagation chain analysis is shown in Figure 2. It is performed on the service call graph in the following three steps.

1) Entry Node Analysis. Treat the initial anomalous service in the service call graph as the entry node. Conduct service anomaly detection for each of its neighboring nodes (i.e., those services that directly invoke or are directly invoked by the initial anomalous service) in terms of different anomaly types and the corresponding quality metrics. For each detected neighboring anomalous node, if its upstream/downstream relationship with the entry node is consistent with the anomaly propagation direction of the detected anomaly type, start an anomaly propagation chain analysis from it and make the detected anomaly type as the anomaly type of the anomaly propagation chain. For example, for the entry node S_5 in Figure 2, two neighboring anomalous nodes S_4 and S_7 are detected and their anomaly types are traffic anomaly and performance anomaly respectively. As S_4 is an upstream node of S_5 and its relationship with S_5 is consistent with the anomaly propagation direction of traffic anomaly (from upstream to downstream), a traffic anomaly propagation chain analysis is started from S_4 . Similarly, a performance anomaly propagation chain analysis is started from S_7 .

2) Anomaly Propagation Chain Extension. For each anomaly propagation chain, iteratively extend it by backtracking along the anomaly propagation direction from the starting point (i.e., a detected neighboring anomalous node of the initial anomalous service). In each iteration, conduct service anomaly detection for all the upstream/downstream nodes of the current node in terms of the anomaly type

of the current anomaly propagation chain; for each detected upstream/downstream anomalous node, add it to the current anomaly propagation chain. The extension of the anomaly propagation chain ends when no more nodes can be added to the chain. For example, for the anomaly propagation chain of S_4 the extension ends with S_1 , which is an upstream node of S_4 and conforms to the propagation direction of traffic anomaly to S_4 . Similarity, for the anomaly propagation chain of S_7 the extension ends with S_9 and S_{10} .

3) Candidate Root Causes Output. When all the anomaly propagation chain analysis for the initial anomalous service ends, report all the services where the extension of an anomaly propagation chain ends as candidate root causes. For example, the candidate root causes reported for the analysis process shown in Figure 2 include S_1 , S_9 , and S_{10} .

B. Service Anomaly Detection

During the anomaly propagation chain analysis process, we need to continuously detect whether a service is an anomalous service in terms of certain quality metrics. For example, for the analysis process shown in Figure 2, we first identify S_4 and S_7 in entry node analysis and then S_7 , S_9 , and S_{10} in anomaly propagation chain extension by service anomaly detection. Given an upstream/downstream service S' of the current service S , we detect possible anomaly with S' in terms of different anomaly types by analyzing the historical data of the corresponding quality metrics (i.e., RT, EC, QPS) of the service call between S' and S . Based on the characteristics of different quality metrics, we choose a different analysis model for each anomaly type. These analysis models are designed based on the characteristics of the fluctuations of the corresponding quality metrics as shown in Figure 3. These data are obtained from the monitoring system of Alibaba.

1) Performance Anomaly: Performance anomaly is detected based on the anomalous increase of response time (RT). The challenge here is how to distinguish anomalous fluctuations from normal fluctuations. Figure 3(a) and Figure 3(d) illustrate the fluctuations of response time in two hours and in one week respectively. It can be seen that there are periodic fluctuations in different periods of a day and different days of a week. If we use intuitive rules like 3-sigma, it is likely that some of these periodic fluctuations will be recognized as performance anomalies. To recognize anomalous fluctuations such as the anomalous point in Figure 3(a), we need to consider not only the quality metrics of the current period but also those of the same period of the previous day and the same day of the previous week. Moreover, in historical data the response time is in normal fluctuations most of the time, and there are only a few anomalous fluctuations.

Based on the characteristics of RT, we choose to use OC-SVM (one class support vector machine) to train a prediction model for anomalous RT fluctuations. OC-SVM [20] uses only the information for the target class (normal RT fluctuations) to learn a classifier that can recognize the samples belonging to the target class and identify the others as outliers (anomalous RT fluctuations). OC-SVM is known to have the advantages

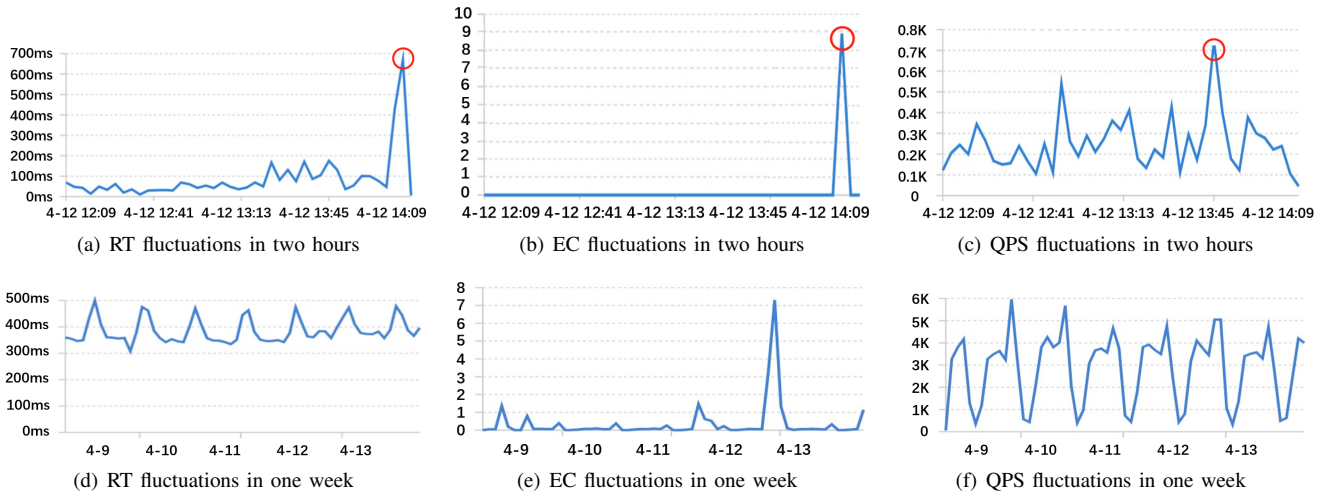


Fig. 3. Fluctuation of Quality Metrics

of simple modelling, strong interpretability, and better generalization ability than clustering methods. We define and use the following four kinds of features for the model.

- **Number of Over-Max Values:** the number of RT values in the current detection window that exceed the maximum RT value of a given comparison period.
- **Delta of Maximum Values:** the delta between the maximum RT value in the current detection window and the maximum RT value of a given comparison period.
- **Number of Over-Average Values:** the number of RT values in the current detection window that exceed the maximum moving average of the RT values of a given comparison period.
- **Ratio of Average Values:** the ratio of the average of the RT values in the current detection window to the maximum moving average of the RT values of a given comparison period.

We consider the last 10 minutes as the current detection window, i.e., considering the 10 RT values of the last 10 minutes. For the comparison periods, we consider the following three settings: the last one hour before the current detection window, the same hour of the previous day, the same hour of the same day of the previous week. Therefore, we have 12 features in total for the OC-SVM model.

We use the Python based machine learning framework scikit-learn [21] to implement the model. To train the model, we collect 100,000 cases of different services from the monitoring infrastructure of Alibaba. Each case includes a series of RT values of the current detection window and different comparison periods. We also collect 600 cases for verifying the trained model. These cases are different from the training cases and the ratio of positive and negative cases is 1:1. The results of the verification are shown in Table II, which provide the recall, precision, F1-measure, FPR (False Positive Rate) of performance anomaly detection with different coverage of the training set. It can be seen that the model can accurately detect

performance anomalies with only 10-20% of the training set, and it can achieve very high accuracy when using 70% of the training set.

TABLE II
ACCURACY OF PERFORMANCE ANOMALY DETECTION

Coverage	Recall	Precision	F1	FPR
10%	0.86	0.73	0.79	0.32
20%	0.81	0.88	0.84	0.12
30%	0.82	0.92	0.87	0.07
40%	0.83	0.92	0.87	0.08
50%	0.84	0.94	0.89	0.06
60%	0.83	0.94	0.88	0.06
70%	0.86	0.95	0.90	0.05
80%	0.87	0.95	0.91	0.05
90%	0.87	0.95	0.91	0.04
100%	0.87	0.96	0.91	0.03

2) *Reliability Anomaly:* Reliability anomaly is detected based on the anomalous increase of error counts (EC). As shown in Figure 3(b) and Figure 3(e), the value of EC is 0 most of the time. In some cases, even when some errors are raised there is no influence on the business and the service may get back to normal in a short time. For example, service calls fail and errors are raised when the circuit break is open, but the service will get back to normal soon after the load decreases and the circuit break is closed. Therefore, we cannot use the performance anomaly detection model to detect reliability anomalies. Binary classification models like Logic Regression (LR) or Random Forest (RF) can be used. However, as only a small number of EC changes are reliability anomalies, LR models are likely to cause overfitting.

Based on the characteristics of EC, we choose to use Random Forest (RF) to train a prediction model for anomalous EC increases. RF uses multiple decision trees for classification. It can effectively combine multiple features and avoid overfitting. We define and use the following five features for the model. Note that some of the features combine other metrics (e.g.,

RT, QPS) together with EC, as anomalous EC increases often correlate with RT and QPS. Similar to performance anomaly detection, we consider the last 10 minutes as the current detection window.

- **Previous Day Delta Outlier Value:** calculate the deltas of the EC values in the last one hour and the EC values in the same hour of the previous day; use the 3-sigma rule to identify possible outliers of the deltas in the current detection window; if exist return the average of the outliers as the feature value, otherwise return 0.
- **Previous Minute Delta Outlier Value:** calculate the deltas of the EC values and the values of the previous minute in the last one hour; use the 3-sigma rule to identify possible outliers of the deltas in the current detection window; if exist return the average of the outliers as the feature value, otherwise return 0.
- **Response Time Over Threshold:** whether the average RT in the current detection window exceeds a predefined threshold (e.g., 50ms).
- **Maximum Error Rate:** the maximum error rate (i.e., EC divided by number of requests) in the current detection window.
- **Correlation with Response Time:** the Pearson correlation coefficient (see Equation 1) between EC and RT in the current detection window.

We use the Python based machine learning framework scikit-learn [21] to implement the model. To train the model, we collect 1,000 labelled cases of different services from the monitoring infrastructure of Alibaba and the ratio of positive and negative samples is 1:3. We also collect 400 cases for verifying the trained model and the ratio of positive and negative samples is 5:3. The results of the verification are shown in Table III, which provide the recall, precision, F1-measure, FPR (False Positive Rate) of reliability anomaly detection with different coverage of the training set. It can be seen that the model can accurately detect reliability anomalies with only 20% of the training set, and it can achieve very high accuracy when using 40% of the training set.

TABLE III
ACCURACY OF RELIABILITY ANOMALY DETECTION

Coverage	Recall	Precision	F1	FPR
10%	0.29	1	0.44	0.70
20%	0.64	1	0.78	0.35
30%	0.64	1	0.78	0.35
40%	0.86	0.97	0.91	0.14
50%	0.88	0.95	0.91	0.12
60%	0.88	0.90	0.89	0.12
70%	0.93	0.90	0.92	0.09
80%	0.95	0.95	0.94	0.07
90%	0.95	0.95	0.95	0.04
100%	0.98	0.93	0.95	0.02

3) *Traffic Anomaly:* Traffic anomaly is detected based on the anomalous fluctuations of queries per second (QPS). As shown in Figure 3(c) and Figure 3(f), QPS complies with the normal distribution in both short term and long term. There-

fore, we choose to use the 3-sigma rule to detect anomalous QPS fluctuations.

Similar to performance and reliability anomaly detection, we consider the last 10 minutes as the current detection window. Based on the QPS values in the last one hour, we use the 3-sigma rule to detect outliers in the current detection window. To further eliminate false positives, we also check the Pearson correlation coefficient (see Equation 1) between the QPS values and the business metric values of the initial anomalous service. Only when the correlation is higher than a predefined threshold (e.g., 0.9) and 3-sigma outliers are identified in the current detection window, a traffic anomaly is reported for the current service.

C. Pruning Strategy

An anomaly propagation chain may have multiple branches and the number of branches can continue to grow during anomaly propagation chain extension. Therefore, a main challenge in anomaly propagation chain analysis lies in the exponential growth of the number of possible anomaly propagation branches. In these branches, some anomalous services and service calls may be irrelevant to the reported availability issue. To tackle the challenge and improve the efficiency of the analysis, we adopt a pruning strategy to eliminate irrelevant anomalous service calls. The pruning strategy is based on the assumption that two successive edges (service calls) in an anomaly propagation chain have similar change trends of the corresponding quality metrics. For example, in Figure 2 the QPS of the edge $S_1 \rightarrow S_4$ should have similar change trend with the QPS of the edge $S_4 \rightarrow S_5$, otherwise the edge $S_1 \rightarrow S_4$ can be pruned. Similarly, the RT of the edge $S_7 \rightarrow S_9$ and the RT of the edge $S_7 \rightarrow S_{10}$ should have similar change trends with the RT of the edge $S_5 \rightarrow S_7$.

Similar to the correlation estimation in candidate root cause ranking (see Section III), we use the Pearson correlation coefficient [19] to measure the similarity of the change trends of the corresponding quality metrics of two successive service calls. For each service call, the service call graph records a quality metric value (e.g., RT, EC, QPS) per minute. To reflect the latest change trends, we only consider the quality metric values in a latest time window (e.g., the last 60 minutes). The values for a specific quality metric of a service call in the last n minutes thus can be represented by a vector X , where X_i ($1 \leq i \leq n$) represents the value of the i th minute. Given two service calls with the quality metric vectors X and Y , their correlation coefficient can be calculated as Equation 1.

The pruning strategy is executed in the anomaly propagation chain extension process in the following way. Before adding an anomalous node to the current anomaly propagation chain, check the correlation coefficient of the change trends of the corresponding quality metric of the current edge (service call) with the adjacent upstream/downstream edge. If the correlation coefficient is lower than a threshold (e.g., 0.7), the current edge will be pruned and the current anomalous node will not be added to the current anomaly propagation chain. For example, for the example shown in Figure 2, we will check

the correlation coefficient of the QPS change trends of the edge $S_1 \rightarrow S_4$ with the edge $S_4 \rightarrow S_5$. If the correlation coefficient is lower than the threshold, we will not add S_1 to the anomaly propagation chain even if it is a QPS anomalous node. Similarly, we will check the correlation coefficients of the RT change trends of the edge $S_7 \rightarrow S_9$ and $S_7 \rightarrow S_{10}$ with the edge $S_5 \rightarrow S_7$.

V. EXPERIMENTAL STUDY

To evaluate the effectiveness and efficiency of MicroHECL, we conduct a series of experimental studies to answer the following research questions.

- **RQ1 (Localization Accuracy):** How accurate is MicroHECL for locating the root causes and anomaly types of availability issues of microservice systems?
- **RQ2 (Localization Efficiency):** How efficient is the root cause localization process of MicroHECL? How well can it scale with the number of services?
- **RQ3 (Effect of Pruning):** How does the pruning strategy influence the accuracy and efficiency of MicroHECL with different similarity thresholds?

A. Experimental Setup

We collect 75 availability issues from 28 subsystems of the e-commerce system of Alibaba. These subsystems contain 265 services on average (min 7, max 1687, median 82). The time of these availability issues ranges from February to June 2020. All of them have been identified and annotated with anomaly types and root causes by operation engineers. Each of them may have multiple anomaly types, and each anomaly type has a single root cause. Among the 75 availability issues, 37 are caused by performance anomaly, 43 are caused by reliability anomaly, and 21 are caused by traffic anomaly. The construction of the service call graph follows the practice in Alibaba, i.e., collecting service calls in the last 30 minutes and service call metrics (i.e., RT, EC, and QPS) in the last 60 minutes.

To evaluate the accuracy and efficiency, we compare MicroHECL with the following two state-of-the-art baseline approaches.

- **MonitorRank** [6]: It detects root causes of anomalies in service-oriented systems and provides a ranked order list of possible root causes. It uses the historical and current time-series metrics of each service along with service call graph to build an unsupervised model for ranking.
- **Microscope** [12]: It locates anomalous services with a ranked list of possible root causes in microservice systems. It builds and uses service instance causality graph, which captures both communicating (service calls) and non-communicating (services co-located in the same machines) service dependencies. It traverses the graph from the front-end service where an anomaly is reported and ranks possible root causes based the similarity of their metrics with those of the front-end service.

We use the top- k hit ratio ($HR@k$) and mean reciprocal rank (MRR) to measure the accuracy of the root case localization.

- $HR@k$ denotes the probability that the top k result list contains the root cause. Let A be the set of availability issues, r_i be the root cause of the i th availability issue, $Rank_i^k$ be the top k result list of the i th availability issue, $HR@k$ can be calculated as follows.

$$HR@k = \frac{1}{|A|} \sum_{i=1}^{|A|} (r_i \in Rank_i^k) \quad (2)$$

- MRR is the multiplicative inverse of the rank of the first correct answer. If the correct answer is not included in the returned list, the rank can be regarded as positive infinity. Let A be the set of availability issues and $Index_i$ be the rank of the root cause in the returned list of the i th availability issue, MRR can be calculated as follows.

$$MRR = \frac{1}{|A|} \sum_{i=1}^{|A|} \frac{1}{Index_i} \quad (3)$$

All the experiments are conducted on a cluster with 15 virtual machines on Alibaba Cloud. Each virtual machine is equipped with 4 Intel Xeon 2.50GHz CPU, 8 GB RAM and 60 GB Disk, and running Alibaba Group Enterprise Linux Server release 6.2. The default settings of the thresholds are the following: the RT threshold used in reliability anomaly detection is 50ms; the correlation threshold used in the traffic anomaly detection is 0.9; the correlation threshold used in the pruning strategy is 0.7.

B. Localization Accuracy (RQ1)

Table IV shows the results of the overall accuracy evaluation of MicroHECL and the two baseline approaches. It can be seen that with MicroHECL the top-1, top-3, and top-5 hit ratios are 0.48, 0.67, and 0.72, respectively, and the MRR is 0.58. MicroHECL significantly outperforms the baseline approaches in terms of all these metrics.

TABLE IV
OVERALL ACCURACY EVALUATION (BEST SCORES ARE IN BOLDFACE)

Metric	MicroHECL	MonitorRank	Microscope
HR@1	0.48	0.32	0.35
HR@3	0.67	0.40	0.49
HR@5	0.72	0.43	0.59
MRR	0.58	0.37	0.44

We also evaluate the accuracy of the three approaches for different anomaly types as shown in Table V. It can be seen that MicroHECL outperforms the two baseline approaches for all the three types of anomalies. The advantages of MicroHECL are the most significant for performance anomaly and the least significant for traffic anomaly.

The advantages of MicroHECL can be explained by its specially designed mechanisms for service anomaly detection and anomaly propagation analysis. MonitorRank considers the anomaly correlations between front-end services and back-end services and the correlations decrease with the propagation. Therefore, it tends to choose services that are closer

TABLE V
ACCURACY EVALUATION BY ANOMALY TYPES (BEST SCORES ARE IN BOLDFACE)

Metric	MicroHECL	MonitorRank	Microscope
Performance Anomaly			
HR@1	0.51	0.27	0.30
HR@3	0.65	0.30	0.46
HR@5	0.70	0.32	0.54
MRR	0.61	0.31	0.39
Reliability Anomaly			
HR@1	0.30	0.26	0.26
HR@3	0.56	0.33	0.44
HR@5	0.70	0.35	0.52
MRR	0.44	0.32	0.35
Traffic Anomaly			
HR@1	0.46	0.32	0.36
HR@3	0.55	0.36	0.46
HR@5	0.59	0.41	0.50
MRR	0.54	0.38	0.46

to the front-end service in the anomaly propagation chains. In contrast, MicroHECL considers the correlations between two successive service calls in the pruning. Microscope uses the 3-sigma rule to detect anomalies and thus may produce false positives when occasional fluctuations occur. In contrast, MicroHECL uses more specific and sophisticated models for different types of anomalies

C. Localization Efficiency (RQ2)

The 75 availability issues are collected from 28 subsystems and these subsystems have different numbers of services. To evaluate the efficiency and scalability of MicroHECL, we analyze the execution time of MicroHECL and the two baseline approaches for the 75 availability issues and investigate how the time changes with the size (service number) of the target system. The results are shown in Figure 4. Note that there may be multiple availability issues collected from the same subsystem and each of them is indicated by a point.

In general, the execution time of MicroHECL is 22.3% less than that of Microscope and 31.7% less than that of MonitorRank. For the subsystems of different sizes MicroHECL uses the least time among the three approaches for most availability issues. The two curves in Figure 4 show the changes of the time differences of MicroHECL with the two baseline approaches respectively. It can be seen that the advantage of MicroHECL is not significant when the number of services is less than 250; when the number of services exceeds 250, the advantage of MicroHECL is getting more and more significant with the increase of service number. Moreover, the execution time of MicroHECL (also the two baseline approaches) increases linearly with the increase of service number, showing a good scalability.

The advantages of MicroHECL can be explained by its specially designed mechanisms for anomaly propagation analysis. MonitorRank traverses the service call graph using random walk algorithm and needs to calculate the correlation scores for a large number of the nodes in the graph. This process is time-

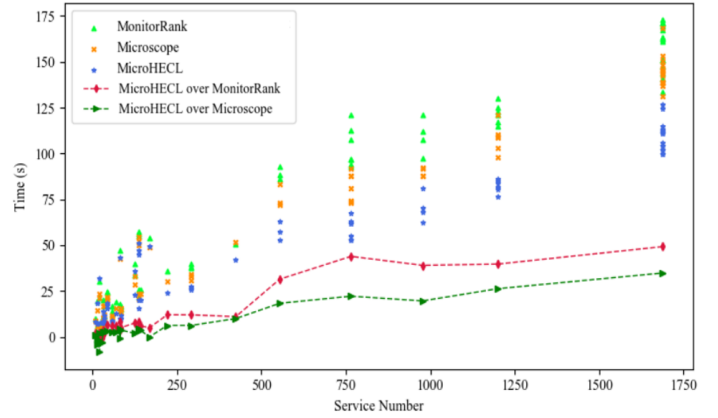


Fig. 4. Detection Time Changes with Num of Nodes

consuming when the graph contains many nodes. Microscope traverses all the neighboring nodes of an anomalous node, both upstream and downstream. Moreover, it has no pruning strategy. In contrast, MicroHECL considers only a single direction (upstream or downstream) for each anomaly type in anomaly propagation chain extension and uses a pruning strategy to eliminate branches that are likely irrelevant to the anomaly propagation.

D. Effect of Pruning (RQ3)

The pruning strategy is a key for the accuracy and efficiency of root cause localization. We evaluate the effect of the pruning strategy by analyzing how the accuracy and time of root cause localization change with the threshold of correlation coefficient. To this end, we run MicroHECL to analyze the 75 availability issues with different threshold settings and measure the accuracy and time. We choose the top-3 hit ratio (i.e., HR@3) as the indicator of accuracy.

The results of the evaluation are shown in Figure 5. It can be seen that both the accuracy and time decrease with the increase of the threshold, as more services and service call edges are pruned in the analysis process and less services are reached and considered. It can also be seen that the accuracy remains 0.67 when the threshold increases from 0 to 0.7, while the time decreases from 75 seconds to 46 seconds. This analysis confirms the effectiveness of the pruning strategy, which can significantly improve the efficiency of root cause localization while keeping the accuracy. And the best threshold is 0.7 for these availability issues.

E. Threats to Validity

A major threat to the internal validity of our studies lies in the implementation of the two baseline approaches. We implement the two approaches by ourselves based on the descriptions in their papers and thus our implementations may not fully embody the approaches. Another threat to the internal validity lies in the interactions with the monitoring infrastructure, i.e., EagleEye. In the experiments, MicroHECL together with the two baseline approaches need to obtain data from EagleEye in real time and there is thus uncertain latency due

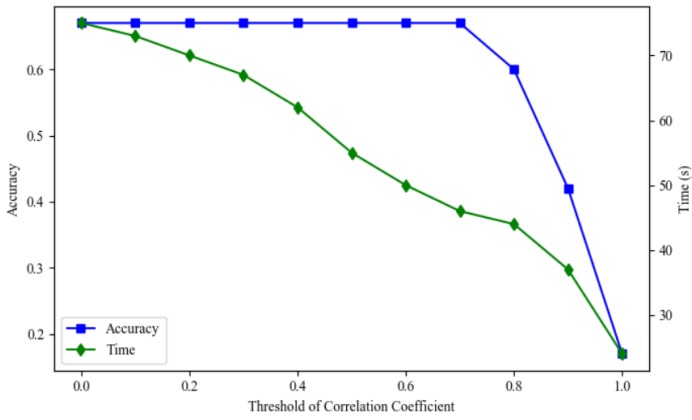


Fig. 5. Evaluation of the Effect of the Pruning Strategy

to limit circuit. This threat applies to all the three approaches equally.

A major threat to the external validity of our studies lies in the fact that all the studies are conducted with the data and availability issues collected from the e-commerce system of Alibaba. Microservice systems of other companies or domains may have different characteristics, for example, in the fluctuations of quality metrics and propagation of service anomalies. The accuracy and efficiency achieved in our studies may not apply to other microservice systems.

VI. PRACTICAL APPLICATION

MicroHECL has been implemented and deployed as a distributed system in Alibaba based on EagleEye and other infrastructures. To support finer-grained localization, the system also considers local middleware components (including databases, message queues, caches) of services as the targets of localization. These components and their interactions with services are incorporated into the service call graph and the quality metrics of the interactions are also recorded for analysis.

MicroHECL has been deployed in Alibaba for more than 5 months and used to handle more than 600 availability issues. For each availability issue the system recommends the top-3 candidate root causes for the developers. Before the application of MicroHECL, the developers need to use the service call graph visualization and quality metrics analysis provided by the web interfaces of EagleEye to manually identify possible root causes. For each availability issue, the typical time of root cause localization and confirmation is over 30 minutes. With MicroHECL, the system can produce root cause recommendations in 76 seconds on average and the developers usually can confirm the results and start fault repairing and recovery in 5 minutes after an availability issue is detected. The hit ratio (i.e., HR@3) of MicroHECL is 68%. Feedback from the developers shows that most of them trust the recommendations of the system and treat the recommendations as the root causes by default.

We analyze the cases for which MicroHECL cannot accurately recommend the root causes and find that the approach

and the system need to be improved from multiple aspects. First, it is possible that availability issues are detected but there are no anomalies with the quality metrics of services. Second, anomalous fluctuations may accumulate slowly and exceed the predefined time window. For example, an anomalous fluctuation may start more than one hour before an availability issue is detected, thus the anomalous fluctuation may not be detected by considering the changes of quality metrics in the last one hour. More advanced techniques are required to improve the current approach.

VII. RELATED WORK

Traditional anomaly detection approaches for distributed and cloud-based systems are mainly based on log analysis. Fu *et al.* [22] proposed an unstructured log analysis technique of anomalies detection for distributed systems. Xu *et al.* [23] introduced LogDC, a log model based problem diagnosis tool for cloud applications with the full-lifecycle Kubernetes logs. Jia *et al.* [24] proposed an approach for automatic anomaly detection based on logs. It raises anomaly alerts on observing deviations from the hybrid model which captures normal execution flows inter and intra services. These approaches build anomaly detection modes with log parsing. They cannot detect the anomalies with quality metrics of services, nor can they support anomaly propagation analysis for root cause localization.

Root cause analysis in practice often relies on visualization of logs and traces. Zhou *et al.* [2] presented an improved visualization analysis approach for fault analysis and debugging of microservice systems. Wang *et al.* [8] demonstrated Grano, an end-to-end anomaly detection and root cause analysis system for cloud native distributed data platform by providing a holistic view of the system component topology, alarms and application events. Guo *et al.* [3] proposed a graph-based trace analysis approach for understanding architecture and diagnosing problems of microservice systems. These approaches help the developers to manually detect the root causes with the aid of trace and log visualization.

Recently, trace analysis based approaches have been proposed for automatic root cause localization for microservice or service-based systems. Pham *et al.* [4] introduced a fault localization approach based on trace similarity comparison. It collects a large number of traces labelled with fault types and root causes through fault injection and conducts fault localization by comparing between faulty traces and historical traces. Gan *et al.* [25] presented Seer, an online cloud performance debugging system for QoS violations by learning spatial and temporal patterns from traces. Zhou *et al.* [5] proposed MEPFL, an approach that can predict latent errors and locate root causes for microservice applications by learning from system trace logs. It predicts latent errors, faulty microservices, and fault types for trace instances captured in the production environment using models trained based on a set of features defined on microservice traces. These approaches often need to collect and analyze a large number of traces to extract anomaly

patterns or train prediction models, thus are not efficient for large-scale microservice systems.

Some researchers have proposed approaches that use service call graphs or causality graphs for automatic root cause localization for microservice or service-based systems. Álvaro *et al.* [9] presented a root cause analysis framework, based on graph representations of service-oriented and microservice architectures. It analyzes anomalies by comparing the similarity between the fault subgraph and the historical subgraph. Li *et al.* [10] proposed MicroRCA, a system to locate root causes of performance issues in microservice systems. MicroRCA infers the root causes by correlating the performance symptoms with the corresponding resource utilization based on an attributed graph. Wang *et al.* [11] proposed CloudRanger, a system dedicated for cloud native systems. CloudRanger identifies the culprit services which are responsible for cloud incidents by constructing the impact graph among services with a dynamic causal relationship analysis approach. Kim *et al.* [6] introduced MonitorRank, an algorithm that can reduce the time, domain knowledge, and human effort required to find the root causes of anomalies in service-oriented architectures. MonitorRank finds root causes with service call graph and historical and current time-series metrics of services. Lin *et al.* [12] proposed Microscope, a system to identify and locate the anomalous services with a ranked list of possible root causes in microservice environments. Microscope infers the causes of performance problems in real time by constructing a service causal graph. The preceding approaches share some similar ideas with ours. However, most of these approaches detect anomalous services by simply analyzing the correlations between metrics or statistical properties and do not consider the characteristics of different anomaly types and quality metrics. Moreover, these approaches need to traverse most of the nodes in the graph, thus may be inefficient for large-scale systems. In contrast, we combine machine learning and statistical methods and design a customized model for the detection of each type of anomalies and adopt a pruning strategy to improve the efficiency of anomaly propagation chain analysis.

VIII. CONCLUSION

In this paper, we propose a high-efficient root cause localization approach for availability issues of microservice systems, called MicroHECL. It dynamically constructs a service call graph and analyzes possible anomaly propagation chains by traversing the graph along anomalous service calls. Different from existing approaches, MicroHECL uses customized models based on machine learning and statistical methods to detect different types of service anomalies (i.e., performance, reliability, traffic) and employs a pruning strategy to eliminate irrelevant service calls in anomaly propagation chain analysis. The accuracy and efficiency of MicroHECL have been confirmed by both experimental studies and the practical application in Alibaba. Future work will be focused on efficiently combing traces, logs, and metrics for more accurate and explainable root cause localization.

REFERENCES

- [1] P. D. Francesco, I. Malavolta, and P. Lago, "Research on architecting microservices: Trends, focus, and potential for industrial adoption," in *ICSA*, 2017, pp. 21–30.
- [2] X. Zhou, X. Peng, T. Xie, J. Sun, C. Ji, W. Li, and D. Ding, "Fault analysis and debugging of microservice systems: Industrial survey, benchmark system, and empirical study," *IEEE Transactions on Software Engineering*, pp. 1–1, 2018.
- [3] X. Guo, X. Peng, H. Wang, W. Li, H. Jiang, D. Ding, T. Xie, and L. Su, "Graph-based trace analysis for microservice architecture understanding and problem diagnosis," in *ESEC/SIGSOFT FSE*, 2020.
- [4] C. Pham, L. Wang, B. Tak, S. Baset, C. Tang, Z. T. Kalbarczyk, and R. K. Iyer, "Failure diagnosis for distributed systems using targeted fault injection," *IEEE Trans. Parallel Distributed Syst.*, vol. 28, no. 2, pp. 503–516, 2017.
- [5] X. Zhou, X. Peng, T. Xie, J. Sun, C. Ji, D. Liu, Q. Xiang, and C. He, "Latent error prediction and fault localization for microservice applications by learning from system trace logs," in *ESEC/SIGSOFT FSE*, 2019, pp. 683–694.
- [6] M. Kim, R. Sumbaly, and S. Shah, "Root cause detection in a service-oriented architecture," in *SIGMETRICS*, 2013, pp. 93–104.
- [7] P. Chen, Y. Qi, P. Zheng, and D. Hou, "Causeinfer: Automatic and distributed performance diagnosis with hierarchical causality graph in large distributed systems," in *INFOCOM*, 2014, pp. 1887–1895.
- [8] H. Wang, P. Nguyen, J. Li, S. Köprü, G. Zhang, S. Katariya, and S. Ben-Romdhane, "GRANO: interactive graph-based root cause analysis for cloud-native distributed data platform," *Proc. VLDB Endow.*, vol. 12, no. 12, pp. 1942–1945, 2019.
- [9] 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," *J. Syst. Softw.*, vol. 159, 2020.
- [10] L. Wu, J. Tordsson, E. Elmroth, and O. Kao, "Microrca: Root cause localization of performance issues in microservices," in *NOMS*, 2020, pp. 1–9.
- [11] P. Wang, J. Xu, M. Ma, W. Lin, D. Pan, Y. Wang, and P. Chen, "Cloudranger: Root cause identification for cloud native systems," in *CCGRID*, 2018, pp. 492–502.
- [12] J. Lin, P. Chen, and Z. Zheng, "Microscope: Pinpoint performance issues with causal graphs in micro-service environments," in *ICSOE*, 2018, pp. 3–20.
- [13] Docker.Com, "Docker," 2020. [Online]. Available: <https://docker.com/>
- [14] Kubernetes.Io, "Kubernetes," 2020. [Online]. Available: <https://kubernetes.io/>
- [15] Istio.Io, "Istio," 2020. [Online]. Available: <https://istio.io/>
- [16] Z. Cai, W. Li, W. Zhu, L. Liu, and B. Yang, "A real-time trace-level root-cause diagnosis system in alibaba datacenters," *IEEE Access*, vol. 7, pp. 142 692–142 702, 2019.
- [17] Nginx.Org, "Nginx," 2020. [Online]. Available: <http://nginx.org/>
- [18] Apache.Org, "Flink," 2020. [Online]. Available: <https://flink.apache.org/>
- [19] H. Abe and S. Tsumoto, "Analyzing behavior of objective rule evaluation indices based on a correlation coefficient," in *KES (2)*, 2008, pp. 758–765.
- [20] R. Vlasveld, "Introduction to one-class support vector machines," 2013. [Online]. Available: <http://rvlasveld.github.io/blog/2013/07/12/introduction-to-one-class-support-vector-machines/>
- [21] Scikit-learn.Org, "Scikitlearn," 2020. [Online]. Available: <https://scikit-learn.org/>
- [22] Q. Fu, J. Lou, Y. Wang, and J. Li, "Execution anomaly detection in distributed systems through unstructured log analysis," in *ICDM*, 2009, pp. 149–158.
- [23] J. Xu, P. Chen, L. Yang, F. Meng, and P. Wang, "Logdc: Problem diagnosis for declaratively-deployed cloud applications with log," in *ICEBE*, 2017, pp. 282–287.
- [24] T. Jia, P. Chen, L. Yang, Y. Li, F. Meng, and J. Xu, "An approach for anomaly diagnosis based on hybrid graph model with logs for distributed services," in *ICWS*, 2017, pp. 25–32.
- [25] Y. Gan, Y. Zhang, K. Hu, D. Cheng, Y. He, M. Pancholi, and C. Delimitrou, "Seer: Leveraging big data to navigate the complexity of performance debugging in cloud microservices," in *ASPLOS*, 2019, pp. 19–33.