# Graph-based Incident Aggregation for Large-Scale Online Service Systems

Zhuangbin Chen*, Jinyang Liu*, Yuxin Su*†, Hongyu Zhang‡
Xuemin Wen¶, Xiao Ling¶, Yongqiang Yang¶, Michael R. Lyu*
*The Chinese University of Hong Kong, Hong Kong, China, {zbchen, jyliu, yxsu, lyu}@cse.cuhk.edu.hk
‡The University of Newcastle, NSW, Australia, hongyu.zhang@newcastle.edu.au
¶Huawei, China, {wenxuemin, lingxiao1, yangyongqiang}@huawei.com

*Abstract*—As online service systems continue to grow in terms of complexity and volume, how service incidents are managed will significantly impact company revenue and user trust. Due to the cascading effect, cloud failures often come with an overwhelming number of incidents from dependent services and devices. To pursue efficient incident management, related incidents should be quickly aggregated to narrow down the problem scope. To this end, in this paper, we propose GRLIA, an incident aggregation framework based on graph representation learning over the cascading graph of cloud failures. A representation vector is learned for each unique type of incident in an unsupervised and unified manner, which is able to simultaneously encode the topological and temporal correlations among incidents. Thus, it can be easily employed for online incident aggregation. In particular, to learn the correlations more accurately, we try to recover the complete scope of failures' cascading impact by leveraging fine-grained system monitoring data, i.e., Key Performance Indicators (KPIs). The proposed framework is evaluated with real-world incident data collected from a large-scale online service system of Huawei Cloud. The experimental results demonstrate that GRLIA is effective and outperforms existing methods. Furthermore, our framework has been successfully deployed in industrial practice.

*Index Terms*—Cloud computing, online service systems, incident management, graph representation learning

## I. Introduction

In recent years, IT enterprises started to deploy their applications as online services on cloud, such as Microsoft Azure, Amazon Web Services, and Google Cloud Platform. These cloud computing platforms have benefited many enterprises by taking over the maintenance of IT and infrastructure and allowing them to improve their core business competence. However, for large-scale online service systems, failures are inevitable, which may lead to performance degradation or service unavailability. Whether or not the service failures are properly managed will have a great impact on the company's revenue and users' trust. For example, an hour episode of downtime in Amazon led to a loss of 100+ million dollars [1].

When a failure happens, system monitors will render a large number of incidents to capture different failure symptoms [2]–[4], which can help engineers quickly obtain a big picture of the failure and pinpoint the root cause. For example, "Special instance cannot be migrated" is a critical network failure in Virtual Private Cloud (VPC) service, and the incident "Tunnel

bearing network pack loss" is a signal for this network failure, which is caused by the breakdown of a physical network card on the tunnel path. Due to the large scale and complexity of online service systems, the number of incidents is overwhelming the existing incident management systems [2], [4], [5]. When a service failure occurs, aggregating related incidents can greatly reduce the number of incidents that need to be investigated. For example, linking incidents that are caused by a hardware issue can provide engineers with a clear picture of the failure, e.g., the type of the hardware error or even the specific malfunctioning components. Without automated incident aggregation, engineers may need to go through each incident to discover the existence of such a problem and collect all related incidents to understand it. Moreover, incident aggregation can also facilitate failure diagnosis. In cloud systems, some trivial incidents are being generated continuously, and multiple (independent) failures can happen at the same time. Identifying correlated incidents can therefore accelerate the process of root cause localization.

To aggregate related incidents, one straightforward way is to measure the text similarity between two incident reports [3], [4]. For example, incidents that share similar titles are likely to be related. Besides textual similarity, system topology (e.g., service dependency, network IP routing) is also an important feature to resort to. Due to the dependencies among online services, failures often have a cascading effect on other inter-dependent services. A service dependency graph can help track related incidents caused by such an effect. However, as cloud systems often possess certain ability of fault tolerance, some services may not report incidents, impeding the tracking of failures' impact (to be explained in Section II). This issue is ubiquitous in production systems, which has not yet been properly addressed in existing work. Moreover, the patterns of incidents are collectively influenced by different factors, such as their topological and temporal locality. Existing work [3], [4] combine them by a simple weighted sum, which may not be able to reveal the latent correlations among incidents.

In this work, we propose GRLIA (stands for Graph Representation Learning-based Incident Aggregation), which is an incident aggregation framework to assist engineers in failure understanding and diagnosis. Different from the existing work of alert storm handling [4] and linked incident identifica-

tion [3], we do not rely on incidents' textual similarity. More-over, we learn incidents' topological and temporal correlations in a unified manner (instead of by a weighted combination). Traditional applications of graph representation learning often learn the semantics of a fixed graph. Unlike them, we propose to learn a feature representation for each unique type of incident, which can appear in multiple places of the graph. The representation encodes the historical co-occurrence of incidents and their topological structure. Thus, they can be naturally used for incident aggregation in online scenarios. To track the impact graph of a failure (i.e., the incidents triggered by the failure), we exploit more fine-grained system signals, i.e., KPIs, as a piece of auxiliary information to discover the scope of its cascading effect. KPIs profile the impact of failures in a more sophisticated way. Therefore, if two services exhibit similar abnormal behaviors (characterized by incidents and KPIs), they should be suffering from the same problems even if no incidents have been reported. Finally, we apply community detection algorithms to find the scope of different failures.

To sum up, this work makes the following major contributions:

- We propose to identify service failures' impact graph, which consists of the incidents that originate from the same failures. Such an impact graph helps us obtain a complete picture of failures' cascading effect. To this end, we combine incidents with KPIs to measure the behavioral similarity between services. Community detection algorithms are then applied to determine the failure-impact graph of different failures automatically.
- We propose GRLIA, an incident aggregation framework based on graph representation learning. The embedding vector for each unique type of incident is learned in an unsupervised and unified fashion, which encodes its interactions with other incidents in temporal and topological dimensions. Online incident aggregation can then be naturally performed by calculating their distance. The implementation of GRLIA is available on GitHub [6].
- We conduct experiments with real-world incidents collected from Huawei Cloud, which is a large-scale cloud service provider. The results demonstrate the effectiveness of the proposed framework. Furthermore, our framework has been successfully incorporated into the incident management system of Huawei Cloud. Feedback from on-site engineers confirms its practical usefulness.

The remainder of this paper is organized as follows. Section II introduces the background and problem statement of this paper. Section III describes the proposed framework. Section IV shows the experiments and experimental results. Section V presents our success story and lessons learned from practice. Section VI discusses the related work. Finally, Section VII concludes this work.

## II. BACKGROUND AND PROBLEM STATEMENT

### A. Topology of Large-scale Online Service Systems

Cloud vendors provide a variety of online services to customers worldwide. In general, there are three main models
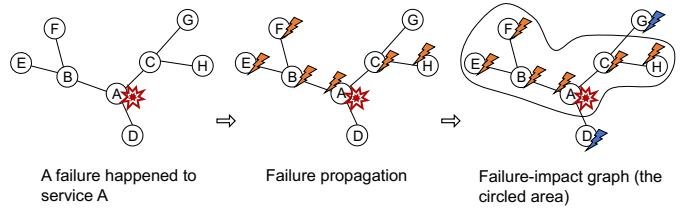


Fig. 1. An illustration of service failures' cascading effect. The irregular circle in the third subfigure shows the failure-impact graph.
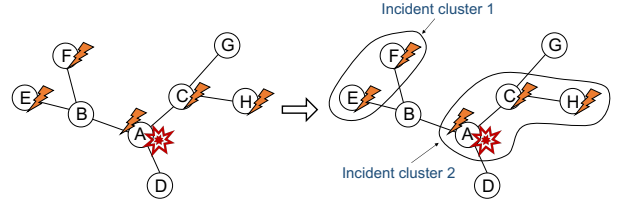


Fig. 2. An example of incomplete failure-impact graph

of cloud-based services, namely, Software as a Service (SaaS), Platform as a Service (PaaS), and Infrastructure as a Service (IaaS) [7]. Online service systems often possess a hierarchical topology, i.e., the stack of application, platform, and infrastructure layers. Each service embodies the integration of code and data required to execute a complete and discrete functionality. For example, in the application layer, the services provided to customers can be a user application, a microservice, or even a function; in the platform layer, the services can be a container or a database; in the infrastructure layer, the services can be a virtual machine or storage. Different services communicate through virtual networks using protocols such as Hypertext Transfer Protocol (HTTP) and Remote Procedure Call (RPC). Such communications among services constitute the complex topology of large-scale online service systems.

### B. Cascading Effect of Service Failures

With such a topology, a failure occurring to one service tends to have a cascading effect across the entire system. Representative service failures include slow response, request timeout, service unavailability, etc., which could be caused by capacity issues, configuration errors, software bugs, hardware faults, etc. To quickly understand failure symptoms, a large number of monitors are configured to monitor the states of different services in a cloud system [2]. A monitor will render an incident when certain predefined conditions (e.g., "CPU utilization rate exceeds 80%") are met. Typical configurations of monitors include setting thresholds for specific metrics (e.g., RPC latency, error counter), checking service/device availability or status, etc. When a failure happens, the monitors often render a large number of incidents. These incidents are triggered by the common root cause and describe the failure from different aspects. Thus, they can be aggregated to help engineers understand and diagnose the failure.

In this paper, we model the set of incidents triggered by a failure as its impact graph (or failure-impact graph), as illus-

TABLE I
EXAMPLES OF INCIDENT AGGREGATION

| No. | | Incident Title | Time | Pod | Severity |
|---|---|---|---|---|---|
| 1 | | Virtual machine is in abnormal state | 2020/10/09 19:40 | pod01 | Low |
| 2 | | Virtual network interface receive lost ratio over 20% | 2020/10/09 19:40 | pod02 | High |
| 3 | Incident group 1 | Traffic burst seen in Nginx node | 2020/10/09 19:40 | pod02 | Low |
| 4 | | Traffic burst seen in LVS (Linux Virtual Server) node | 2020/10/09 19:41 | pod09 | Medium |
| 5 | | OSPF (Open Shortest Path First) protocol state change | 2020/10/09 19:41 | pod04 | Medium |
| 6 | | Excessive I/O delay of storage disk | 2020/10/12 14:34 | pod09 | Medium |
| 7 | | Component failure | 2020/10/12 14:34 | pod05 | High |
| 8 | Incident group 2 | Hard disk failure | 2020/10/12 14:34 | pod09 | Medium |
| 9 | | Database account login error | 2020/10/12 14:34 | pod18 | Medium |
| 10 | | Monitor detected customer-impacting incidents for Storage in [AZ1] | 2020/10/12 14:35 | pod10 | Medium |

trated by Fig. 1. Specifically, service $A$ encounters a failure, and the impact propagates to other services along the system topology. The circled area indicates the impact graph of the failure, where irrelevant incidents (in blue) in service $D$ and $G$ are excluded. In general, the system topology can have many different forms, such as the dependencies of services [8], the configured IP routing of a cloud network [9], etc. Intuitively, it might seem that the impact graph can be easily constructed by tracing incidents along the system topology. However, our industrial practices reveal that they are usually incomplete. An example is given in Fig. 2, where service $B$ occasionally fails to report any incident during the failure. Existing approaches may perceive it as two separate failures, which is undesirable. We have summarized the following two main reasons for the missing incidents:

- System monitors that report incidents are configured with rules predefined by engineers. Due to the diversity of cloud services and user behaviors, the impact of a failure may not meet the rules of some monitors. For example, if an application triggers an incident when its CPU usage rate exceeds 80%, then any value below the threshold will be unqualified. As a consequence, the monitors will not report any incident, and thus, the tracking of the failure's impact is blocked.
- To ensure the continuity of online services, cloud systems are designed to have a certain fault tolerance capability. In this case, service systems can bear some abnormal conditions, and thus, no incidents will be reported. Therefore, the impact of a failure may not manifest itself completely over the system topology.

Recent studies on cloud incident management [2], [10] have demonstrated the incompleteness and imperfection of monitor design and distribution in online service systems. Thus, along the service dependency chain, some services in the middle may remain silent (i.e., report no incident), which impedes the tracking of failure's cascading effect. Therefore, although online service systems generate many incidents, they are often scattered.

## C. Problem Statement

This work aims to assist engineers in failure understanding and diagnosis with online incident aggregation, which is to aggregate incidents caused by the common failure. When services encounter failures, incidents that capture different failure symptoms constitute an essential source for engineers to conduct a diagnosis. However, it is time-consuming and tedious for engineers to manually examine each incident for failure investigation when faced with such an overwhelming number of incidents. Online incident aggregation is to cluster relevant incidents when they come in a streaming manner (i.e., continuously reported by the system). Examples are presented in Table I, where items in blue and gray belong to two groups of aggregated incidents. Particularly, the first group shows a virtual network failure. Note that only the No.3 and No.4 incidents share some words in common, while the others do not. Meanwhile, the second group describes a hardware failure, and more specifically, a storage disk error. Engineers can benefit from such incident aggregation as the problem scope is narrowed down to each incident cluster.

However, accurately aggregating incidents for online service systems is challenging. We have identified three main reasons.

**Background noise**. Although related incidents are indeed generated around the same time, many other cloud components are also constantly rendering incidents. These incidents are mostly trivial issues and therefore become background noise. Incident aggregation based on temporal similarity would suffer from a high rate of false positives.

**Dissimilar textual description**. Text (e.g., incident title and summary) similarity is an essential metric for incident correlation, which has been widely used in existing work [3], [4]. However, in reality, related incidents, especially the critical ones, do not necessarily have similar titles. Failing to correlate such critical incidents greatly hinders root cause diagnosis.

**Unclear failure-impact graph**. To correlate incidents accurately, we need to estimate the impact graph of service failures. As discussed in Section II-B, this task is challenging. Incidents alone are insufficient to completely reflect the impact

failure-impact graph 1    failure-impact graph 2

Service failure detection    Failure-impact graph completion    Graph representation learning    Online incident aggregation
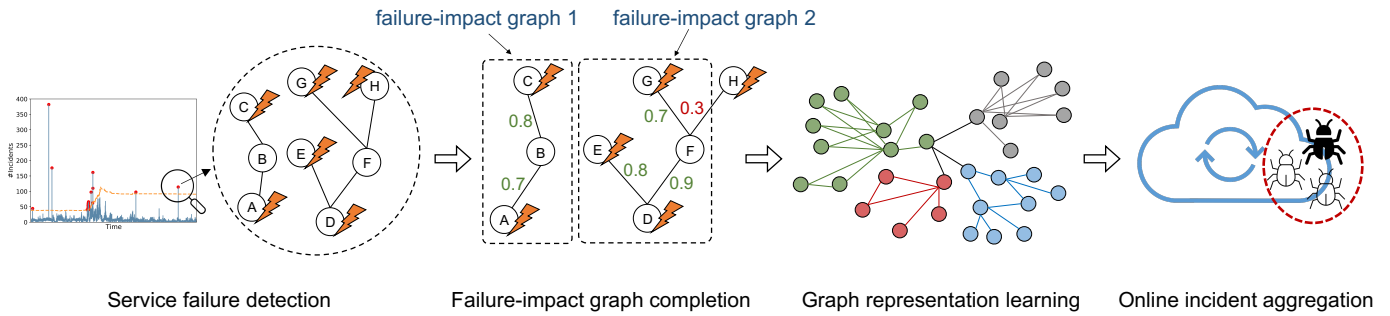
Fig. 3. The overall framework of GRLIA

of failures on the entire system. Therefore, we need to utilize more fine-grained information of the failures.

## III. METHODOLOGY

### A. Overview

In cloud systems, a large number of monitors are configured to continuously monitor the states of its services from different aspects. Many incidents rendered by the monitors tend to co-occur due to their underlying dependencies. For example, some failure symptoms often appear together, and some incidents may develop causal relationship. Our main idea is to capture the co-occurrences among incidents by learning from historical failures. In online scenarios, such correlations can be leveraged to distinguish correlated incidents that are generated in streams.

The overall framework of GRLIA is illustrated in Fig. 3, which consists of four phases, i.e., *service failure detection*, *failure-impact graph completion*, *graph representation learning*, and *online incident aggregation*. The first phase tries to identify the occurrence of service failures and retrieves different types of monitoring data, including incidents, KPI time series, and service system topology. In the second phase, we try to identify the incidents that are triggered by each individual failure detected above. More often than not, it is hard to precisely identify the impact scope of failures (as discussed in Section II-B), which hinders the learning of incidents' correlations. Therefore, we utilize the trends observed in KPI curves to auto-complete the failure-impact graphs. After obtaining the set of incidents associated with each failure, in the third phase, an embedding vector is learned for different types of incidents by leveraging existing graph representation learning models [11], [12]. Such representation encodes not only the temporal locality of incidents, but also their topological relationship. In the final phase, the learned incident representation will be employed for online incident aggregation by considering their cosine similarity and topological distance. In particular, we do not explicitly consider the dynamic change of a system topology because the changes often happen to a small area of the topology, e.g., container creation or kill. GRLIA essentially learns the correlations among incidents, which are also applicable to the changed portion of the topology. Nevertheless, when the system topol-

ogy goes through a significant alteration, our framework is efficient enough to support quick model retraining.

### B. Service Failure Detection

Due to the cascading effect, when service failures occur, a large number of incidents are often reported in a short period of time. Thus, setting a fixed threshold for the average number of reported incidents (e.g., #incidents/min>50) could be a reasonable criterion to detect failures. However, such a design suffers from a trade-off between false positives and false negatives due to online service systems' complex and ever-changing nature [4]. For example, different services have distinct sensitivity to the number of incidents, and continuous system evolution/feature upgrades could change the threshold. Thus, a self-adaptive algorithm is more desirable.

For time-series data, anomalies often manifest themselves as having a large magnitude of upward/downward changes. Extreme Value Theory (EVT) [13] is a popular statistical tool to identify data points with extreme deviations from the median of a probability distribution. It has been applied to predicting unusual events, e.g., severe floods and tornado outbreaks [14], by finding the law of extreme values that usually reside at the tail of a distribution. Moreover, it requires no hand-set thresholds and makes no assumptions on data distribution. In this work, we follow [4], [13] to detect bursts in time series of the number of incidents per minute. As a typical time series anomaly detection problem, other approaches (e.g., [15], [16]) in this field are also applicable. The bursts are regarded as the occurrence of service failures. This algorithm can automatically learn the normality of the data in a dynamic environment and adapt the detection method accordingly. Fig. 3 (phase one) presents an example of service failure detection, where all abnormal spikes are successfully found by the decision boundary (the orange dashed line). For consecutive bins that are marked as anomalies, we regard them as one failure because failures may last for more than one minute. The next phase will distinguish multiple (independent) failures that happen simultaneously. Particularly, the detection algorithm is only required to have a high recall, and the precision is of less importance. It is because the goal of the follow-up two phases is to find the correlations between incidents. Such correlation rules will not be violated even incidents are not appearing together during actual cloud failures.

## C. Failure-Impact Graph Identification

In the first phase, the number of incidents per minute is calculated, and incident bursts are regarded as the occurrence of service failures. For each failure, the incidents collected from the entire system are not necessarily related to it. This is because: 1) while some services are suffering from the failure, others may continuously report incidents (could be trivial and unrelated issues); and 2) multiple service failures could happen simultaneously. Therefore, we need to identify the set of incidents for each individual failure that is generated due to the cascading effect.

To this end, the concept of community detection is exploited. Community detection algorithms aim to group the vertices of a graph into distinct sets, or communities, such that there exist dense connections within a community and sparse connections between communities. Each community represents a collection of incidents rendered by the common service failure, in which the correlations among incidents can be explored. A comparative review of different community detection algorithms is available in [17]. In this work, we employ the well-known *Louvain* algorithm [18], which is based upon modularity maximization. The modularity of a graph partition measures the density of links inside communities compared to links between communities. For weighted graphs, the modularity can be calculated as follows [18]:

$$M = \frac{1}{2m} \sum_{i,j} [W_{i,j} - \frac{k_i k_j}{2m}] \delta(c_i, c_j) \tag{1}$$

where $W_{ij}$ is the weight of the link between node $i$ and $j$, $k_i = \sum_j W_{ij}$ sums the weights of the links associated with node $i$, $c_i$ is the community to which node $i$ is assigned to, $m = \frac{1}{2} \sum_{ij} W_{ij}$, and the $\delta(u,v) = 1$ if $u = v$ and 0 otherwise.

To better understand the identification of failure-impact graph using community detection, an illustrating example is depicted in Fig. 3 (phase two). In this case, except for nodes $B$ and $F$, other nodes all report incidents. By conducting community detection, we obtain two communities: $\{A, B, C\}$ and $\{C, E, F, G\}$, which are regarded as the complete impact graph of their respective failure. The weight between nodes is provided with their link. We can see that intra-community links all have a relatively large weight. Such partition can achieve the best modularity score for this example. Particularly, node $H$ is excluded from the second community due to the small weight of its connection to node $F$.

To apply community detection, the weight between two nodes should be defined. Inspired by [19], we combine KPIs with incidents to calculate the behavioral similarity between two nodes and use the similarity value as the weight. Specifically, the weight is composed of two parts, i.e., incident similarity and KPI trend similarity.

*1) Incident similarity:* Incident similarity is to compare the incidents reported by two nodes. Typically, if two nodes encounter similar errors, they will render similar types of incidents. Jaccard index is employed to quantify such similarity,
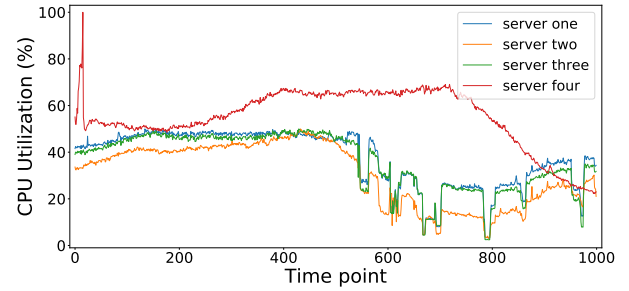


Fig. 4. CPU usage curve of four servers

which is defined as the size of the intersection divided by the size of the union of two incident sets:

$$Jaccard(i,j) = \frac{|inc(i) \cap inc(j)|}{|inc(i) \cup inc(j)|} \tag{2}$$

where $inc(i)$ is the incidents reported by node $i$. In particular, we allow duplicate types of incidents in each set by assigning them a unique number. This is because the distribution of incident types also characterizes the failure symptoms.

*2) KPI trend similarity:* As discussed in Section II, some services may remain silent when failures happen, hindering the tracking of related incidents. To bridge this gap, we resort to KPIs, which are more sophisticated monitoring signals. Intuitively, the KPI trend similarity measures the underlying consistency of cloud components' abnormal behaviors, which cannot be captured by incidents alone. An example is shown in Fig. 4, which records the CPU utilization of four servers. Clearly, the curve of the first three servers exhibits a highly similar trend, while such a trend cannot be observed in server four. The implication is that the first three servers are likely to be suffering from the same issue and thus should belong to the same community. We adopt dynamic time warping (DTW) [20] to measure the similarity between two temporal sequences with varying speeds. We observe the issue of temporal drift between two time series, which is common as different cloud components may not be affected by a failure simultaneously during its propagation. Therefore, DTW fits our scenario.

The remaining problem is which KPIs should be utilized for similarity evaluation. Normal KPIs which record the system's normal status should be excluded as they provide trivial and noisy information. Therefore, EVT introduced in phase one is utilized again to detect anomalies for each KPI. Only the abnormal KPIs shared by two connected cloud components will be compared. Particularly, when there exists more than one type of abnormal KPIs, we use the average similarity score calculated as follows:

$$DTW(i,j) = \frac{1}{K} \sum_{k=1}^{K} dtw(t_k^i, t_k^j) \tag{3}$$

where $K$ is the number of KPIs to compare for node $i$ and $j$, $t_k^i$ is the $k^{th}$ KPI of node $i$, and $dtw(u,v)$ measures the

DTW similarity between two KPI time series $u$ and $v$, which is normalized for path length. The weight $W_{ij}$ between node $i$ and $j$ is computed by taking the weighted sum of the two types of similarities as follows:

$$W_{ij} = \alpha \times Jaccard(i,j) + (1-\alpha) \times DTW(i,j) \quad (4)$$

where the balance weight $\alpha$ is a hyper-parameter. In our experiments, if two nodes both report incidents, we set it as 0.5; otherwise, it is set to be 0, i.e., only the KPI trend similarity is considered.

Finally, for each discovered community, the incidents inside it form the complete impact graph of the service failure. Note that in online scenarios, we cannot directly adopt the techniques introduced in this phase for incident aggregation. This is because they involve a comparison between different KPIs, which are not complete until the failures fully manifest themselves. Thus, the comparison is often delayed and inefficient. Moreover, they can be error-prone without fully considering the historical cases.

### D. Graph-based Incident Representation Learning

After obtaining the impact graph for each service failure (i.e., the actual incidents triggered by it), we can learn the correlations among incidents. Such correlations describe the sets of incidents that tend to appear together. FP-Growth proposed by Han et al. [21] is a standard algorithm to mine such frequent item sets. However, our analysis reveals the following drawbacks it possesses for our problem:

- It is vulnerable to background noise. In production environments, some simple incidents are constantly being reported, e.g., "High CPU utilization rate". These incidents will appear in many transactions (a collection of items that appear together) for FP-Growth. As a result, unrelated incidents might be put into the same frequent item set due to sharing such incidents. These simple incidents cannot be trivially removed as they provide necessary information about a system, and a burst of such incidents can also indicate serious problems.
- It cannot handle incidents with a low frequency. FP-Growth has a parameter called *support*, which describes how frequently an item set is in the dataset. Incident sets with a low support value will be excluded to guarantee the statistical significance of the results. However, more often than not, such incident sets are more important, as they report some critical failures that do not happen frequently.

In online service systems, different resources (e.g., microservices and devices) are naturally structured in graphical forms, such as service dependency and network IP routing. Therefore, graph representation learning [11] can be an ideal solution to deal with the above issues. Graph representation learning is an essential and ubiquitous task with applications ranging from drug design to friendship recommendation in social networks. It aims to find a representation for graph structure that preserves the semantics of the graph. A typical graph representation learning algorithm learns an embedding vector for all nodes of a graph. For example, Chen et al. [3] employed *node2vec* [22] to learn a feature representation for cloud components. Different from them, we propose to learn a representation for each unique type of incident, which can appear in multiple places of the graph. In our framework, we employ *DeepWalk* [23] because of its simplicity and superior performance. DeepWalk belongs to the class of shallow embedding approaches that learn the node embeddings based on random walk statistics. The basic idea is to learn an embedding $\vartheta_i$ for node $v_i$ in graph $\mathcal{G}$ such that:

$$EMB(\vartheta_i, \vartheta_j) \triangleq \frac{e^{\vartheta_i \cdot \vartheta_j}}{\sum_{v_k \in \mathcal{V}} e^{\vartheta_i \cdot \vartheta_k}} \approx p_{\mathcal{G},T}(v_j | v_i) \quad (5)$$

where $\mathcal{V}$ is the set of nodes in the graph and $p_{\mathcal{G},T}(v_j|v_i)$ is the probability of visiting $v_j$ within $T$ hops of distance starting at $v_i$. The loss function to maximize such probability is:

$$\mathcal{L} = \sum_{(v_i, v_j) \in \mathcal{D}} -log(EMB(\vartheta_i, \vartheta_j)) \quad (6)$$

where $\mathcal{D}$ is the training data generated by sampling random walks starting from each node. Readers are referred to the original paper [23] for more details.

For each failure-impact graph, incident sequences are generated through random walk starting from every node inside. In reality, each node usually generates more than one incident when failures happen. Our tailored random walk strategy therefore contains two hierarchical steps. In the first step, a node is chosen by performing random walks on node level; in the second step, an incident will be randomly selected from those reported by the chosen node. Duplicate types of incidents in a node will be kept because frequency is also an important feature of incidents (it impacts the probability of being selected). Following the original setting of [22], we set the walk length as 40, i.e., each incident sequence will contain 40 samples. Finally, the incident sequences will be fed into a Word2Vec model [24] for embedding vector learning. The Word2Vec model has two important hyper-parameters: the window size and the dimension of the embedding vector. We set the window size as ten by following [22] and set the dimension as 128. In particular, by considering the topological distance between incidents, we can alleviate the problem of background noise. This is because as the distance increases, the impact of noisy incidents gradually weakens, while in FP-Growth, all incidents play an equivalent role in a transaction.

### E. Online Incident Aggregation

With the learned incident representation from the last phase, we can conduct incident aggregation in production environments, where the incidents come in a streaming manner. Each group of aggregated incidents represents a specific type of service issue, such as hardware issue, network traffic issue, network interface down, etc. The EVT-based method also plays a role in this phase by continuously monitoring the number of incidents per minute. If it alerts a failure, the online incident

aggregation will be triggered. When two incidents, say $i$ and $j$, appear consecutively, GRLIA measures their similarity. If the similarity score is greater than a predefined threshold, they will be grouped together immediately. In particular, the similarity score consists of two parts, i.e., *historical closeness* (HC) and *topological rescaling* (TR), which are defined as follows:

$$HC(i,j) = \frac{\vartheta_i \cdot \vartheta_j}{\|\vartheta_i\| \times \|\vartheta_j\|}$$
$$TR(i,j) = \frac{1}{max(1, d(i,j) - \mathcal{T})} \quad (7)$$

where $\vartheta_i$ and $\vartheta_j$ are the embedding vectors of incident $i$ and $j$ (as described in Section III-D), respectively; $d(i,j)$ is the topological distance between $i$ and $j$, which is the number of hops along their shortest path in the system topology; and $\mathcal{T}$ is the threshold for considering the penalty of long distance. That is, the topological rescaling becomes effective (i.e., $<1$) only if their distance is larger than $\mathcal{T}$. In our experiments, $\mathcal{T}$ is set as four. Incorrect correlations will be learned if $\mathcal{T}$ is too large, while important correlations will be missed if $\mathcal{T}$ is too small. Our experiments show similar results when $\mathcal{T}$ is in $[3, 6]$. Cosine similarity is adopted to calculate the historical closeness, which is related to their co-occurrences in the past. Finally, the similarity between $i$ and $j$ can be obtained by taking the product of $TR(i,j)$ and $HC(i,j)$:

$$sim(i,j) = TR(i,j) \times HC(i,j)$$
$$= \frac{1}{max(1, d(i,j) - \mathcal{T})} \times \frac{\vartheta_i \cdot \vartheta_j}{\|\vartheta_i\| \times \|\vartheta_j\|} \quad (8)$$

We set an aggregation threshold $\lambda$ for $sim(i,j)$ to consider whether or not two incidents are correlated:

$$cor(i,j) = \begin{cases} 1, & if\ sim(i,j) \geq \lambda \\ 0, & otherwise \end{cases} \quad (9)$$

In our experiments, $\lambda$ is empirically set as 0.7. In particular, the distance of an incident to a group of incidents is defined as the largest value obtained through element-wise comparison.

## IV. EXPERIMENTS

In this section, we evaluate our framework using real-world incidents collected from industry. Particularly, we aim at answering the following research questions.

**RQ1**: How effective is the service failure detection module of GRLIA?

**RQ2**: How effective is GRLIA in incident aggregation?

**RQ3**: Can the failure-impact graph help incident aggregation?

### A. Experiment Setting

*1) Dataset:* Incident aggregation is a typical problem across different online service systems. In this experiment, we select a representative, large-scale system, i.e., the Networking service of Huawei Cloud, to evaluate the proposed framework. Besides offering traditional services such as Virtual Network, VPN

Gateway, it also features intelligent IP networks and other next-generation network solutions. In particular, the service system comprises a large and complex topological structure. In the layer of infrastructure, platform, and software, it has multiple instances of virtual machines, containers, and applications, respectively. In each layer, their dependencies form a topology graph. The cross-layer topology is mainly constructed by their placement relationships, i.e., the mappings between applications, containers, and virtual machines. Like other cloud enterprises, Huawei Cloud's resources are hosted in multiple regions and endpoints worldwide. Each region is composed of several availability zones (isolated locations within regions from which public online services originate and operate) for service reliability assurance. The incident management of the Networking service is also conducted in such a multi-region way, with each region having relatively isolated issues. In this paper, we collect incidents generated between May 2020 and November 2020, during which the Networking service reported a large number of incidents. Although we conduct the evaluation on a single online service system, we believe GRLIA can be easily applied to other online service systems and bring them benefits.

To evaluate the effectiveness of GRLIA, experienced domain engineers manually labeled related incidents. Thanks to the well-designed incident management system with user-friendly interfaces, the engineers can quickly perform the labeling. Note that the manual labels are only required for evaluating the effectiveness of our framework, which is unsupervised. To calculate the KPI trend similarity, we adopt the following KPIs, which are suggested by the engineers:

- *CPU utilization* refers to the amount of processing resources used.
- *Round-trip delay* records the amount of time it takes to send a data packet plus the time it takes to receive an acknowledgement of that data packet.
- *Port in-bound/out-bound traffic rate* refers to the average amount of data coming-in to/going-out of a port.
- *In-bound packet error rate* calculates the error rate of the packet that a network interface receives.
- *Out-bound packet loss rate* calculates the loss rate of the packet that a network interface sends.

These KPIs are representative that characterize the basic states of the Networking service system. In particular, CPU utilization is monitored for different containers and virtual machines, while the remaining KPIs are monitored for the virtual interfaces of each network device. Each KPI is calculated or sampled every minute. We collect two hours of data to measure the KPI trend similarity. Note that the set of KPIs can be tailored for different systems. For example, a database service may also care about the number of failed database connection attempts, the number of SQL queries, etc.

In this paper, we select the largest ten availability zones for experiments, each of which contains a large system topology. Six months of production incidents are collected from the Networking service of Huawei Cloud. The number of distinct

| Dataset | Training period | Evaluation month | #Incidents | #Failures |
|---------|-----------------|------------------|------------|-----------|
| Dataset1 | 2020 May - 2020 July | 2020 Aug. | ~18k/~8k | 105/46 |
| Dataset2 | 2020 May - 2020 Aug. | 2020 Sept. | ~26k/~10k | 151/52 |
| Dataset3 | 2020 May - 2020 Sept. | 2020 Oct. | ~36k/~8k | 203/38 |

incident types is more than 3,000. Similar to [4], [15], [25], we conduct three groups of experiments using incidents reported in the first four months, the first five months, and all months, respectively. In all periods, incident aggregation is applied to the failures that happened in the last month based on the incident representations learned from previous months. Table II summarizes the dataset. For column *#Incidents* (resp. *#Failures*), the first figure calculates the incidents (resp. failures) captured during the training period, while the second figure shows that of the evaluation month. Particularly, some failures are of small scale and can be quickly mitigated, while others are cross-region and become an expensive drain on the company's revenue. We can see each failure is associated with roughly 200 incidents, demonstrating a strong need for incident aggregation.

*2) Evaluation Metrics:* For RQ1, which is a binary classification problem, we employ *precision*, *recall*, and *F1 score* for evaluation. Specifically, precision measures the percentage of incident bursts that are successfully identified as service failures over all the incident bursts that are predicted as failures: $precision = \frac{TP}{TP+FP}$. Recall calculates the portion of service failures that are successfully identified by GRLIA over all the actual service failures: $recall = \frac{TP}{TP+FN}$. Finally, F1 score is the harmonic mean of precision and recall: $F1\ score = \frac{2 \times precision \times recall}{precision+recall}$. $TP$ is the number of service failures that are correctly discovered by GRLIA; $FP$ is the number of trivial incident bursts (i.e., no failure is actually happening) that are wrongly predicted as service failures by GRLIA; $FN$ is the number of service failures that GRLIA fails to discover.

For RQ2 and RQ3, we choose Normalized Mutual Information (NMI) [26], which is a widely used metric for evaluating the quality of clustering algorithms. The value of NMI ranges from 0 to 1 with 0 indicating the worst result (no mutual information) and 1 the best (perfect correlation): $NMI(\Omega, \mathbb{C}) = \frac{2 \times I(\Omega;\mathbb{C})}{H(\Omega)+H(\mathbb{C})}$, where $\Omega$ is the set of clusters, $\mathbb{C}$ is the set of classes, $H(\cdot)$ is the entropy, and $I(\Omega;\mathbb{C})$ calculates and mutual information between $\Omega$ and $\mathbb{C}$.

*3) Implementation:* Our framework is implemented in Python. We parallelize our experiments by assigning availability zones to different processors. The output of each processor is a list of incident sequences generated through random walk, which we merge and feed to a Word2Vec model implemented with Gensim [27], an open-source library for topic modeling and natural language processing. We run our experiments on a machine with 20 Intel(R) Xeon(R) Gold 6148 CPU @ 2.60GHz, and 256GB of RAM. The results show that each phase of our framework takes only a few seconds. The last

phase can even produce results in a real-time manner as it only involves simple vector calculation. Thus, our framework can quickly respond in online scenarios. This demonstrates that GRLIA is of high efficiency.

### B. Comparative Methods

The following methods are selected for comparative evaluation of GRLIA.

- *FP-Growth* [21]. FP-Growth is a widely-used algorithm for association pattern mining. It is utilized as an analytical process that finds a set of items that frequently co-occur in datasets. In our experiments, each impact graph is regarded as a transaction for this algorithm. Given a set of impact graphs, it searches incidents that often appear together, regardless of their distance.

- *UHAS* [4]. This approach is proposed by Zhao et al. aiming at handling alert storms for online service systems. Similar to incident bursts, alert storms also serve as a signal for service failures. Particularly, UHAS employs DBSCAN for alert clustering based on their textual and topological similarity. The textual similarity between two alerts is measured by Jaccard distance. The topological similarity considers two types of topologies, i.e., software topology (service) and hardware topology (server). The topological distance is computed by the shortest path length between two nodes. Finally, a weighted combination of the two types of similarities yields the final similarity score.

- *LiDAR [3].* LiDAR is a supervised method proposed by Chen et al. to identify linked incidents in large-scale online service systems. Specifically, LiDAR is composed of two modules, i.e., textual encoding module and component embedding module. The first module produces similar representations for the text description of linked incidents, which are labeled by engineers. In the evaluation stage, the textual similarity between two incidents is measured by the cosine distance of their representations. The second module learns a representation for the system topology (instead of incidents). The final similarity is calculated by taking a weighted sum of both parts. As LiDAR is supervised, it would be unfair to compare it with other unsupervised methods. Considering the success of Word2Vec model [24], [28] in identifying semantically similar words (in an unsupervised manner), we alter LiDAR to be unsupervised to fit our scenario by representing the text of incidents with off-the-shelf word vectors [29].

### C. Experimental Results

*1)* **RQ1:** *The Effectiveness of GRLIA's Service Failure Detection:* To answer this research question, we compare GRLIA with the fixed thresholding method on three datasets and report precision, recall, and F1 score. Thresholding remains an effective way for anomaly detection in production systems and serves as a baseline in many existing work. Since both methods require no parameter training, we use them to detect failures

| Dataset | Metric | Thresholding | GRLIA |
|---------|--------|--------------|-------|
| Dataset1 | Precision | 0.711 | **0.917** |
| | Recall | 0.913 | **0.957** |
| | F1 Score | 0.799 | **0.937** |
| Dataset2 | Precision | 0.831 | **0.944** |
| | Recall | 0.942 | **0.981** |
| | F1 Score | 0.883 | **0.962** |
| Dataset3 | Precision | 0.648 | **0.925** |
| | Recall | 0.921 | **0.974** |
| | F1 Score | 0.761 | **0.949** |

| Method | Dataset1 | Dataset2 | Dataset3 |
|--------|----------|----------|----------|
| FP-Growth | 0.481 | 0.523 | 0.546 |
| UHAS | 0.697 | 0.71 | 0.707 |
| LiDAR | 0.742 | 0.758 | 0.826 |
| GRLIA | **0.831** | **0.866** | **0.912** |

for both the training data and evaluation data. Particularly, the threshold of the baseline method is #incidents/min>50, which is recommended by field engineers. Moreover, the ground truth is obtained directly from the historical failure tickets, which are stored in the incident management system.

The results are shown in Table III, where GRLIA outperforms the fixed thresholding in all datasets and metrics. In particular, GRLIA achieves an F1 score of more than 0.93 in different datasets, demonstrating its effectiveness in service failure detection. Indeed, we observe that some failures may not always incur a large number of incidents at the beginning. However, if ignored, they could become worse and end up yielding more severe impacts across multiple services. Fixed thresholding does not possess the merit of threshold adaptation based on the context and thus produces many false positives. GRLIA outperforms it for being able to adjust the threshold automatically.

*2) RQ2: The Effectiveness of GRLIA in Incident Aggregation:* We compare the performance of GRLIA against a series of baseline methods for incident aggregation. Table IV shows the NMI values of different experiments. From dataset 1 to 3, GRLIA achieves an NMI score of 0.831, 0.866, and 0.912, respectively, while the best results from the baseline methods are 0.742, 0.758, and 0.826, all attained by LiDAR. LiDAR outperforms UHAS by explicitly considering the entire system topology. Except for UHAS, all approaches achieve better performance with more training data available. This is because UHAS directly works on alert storms when failures are detected. Without learning from the history, it cannot handle complicated scenarios. Recall that both UHAS and LiDAR rely on the textual similarity between incidents. However, in our system, related incidents do not necessarily possess similar text descriptions. For example, there is a clear correlation between the incident "Traffic drops sharply in vRouter" and "OS network ping abnormal" in VPC service, which tends to be missed by them. Moreover, monitors that render incidents are configured by multiple service teams, which further damages the credit of textual similarity. This is particularly true for some critical incidents because they are often tailored for special system errors, which may not be shared across different services. On the other hand, although

GRLIA does not explicitly leverage incident's textual features, our experiments show that it is capable of correlating incidents that share some common words, e.g., "VPC service tomcat port does not exist" and "VPC service tomcat status is dead". This is because such a relationship is reflected in their temporal and topological locality, which can be precisely captured by incidents' representation vectors.

Another observation is that FP-Growth does not fit the task of incident aggregation, whose best NMI score is 0.546. As discussed in Section III-D, this method is not robust against background noise. Indeed, in the system, some trivial incidents (e.g., "Virtual machine is in abnormal state") are continuously being reported, which may connect incidents from distinct groups. Furthermore, many essential incidents are excluded by this method due to low frequency, which is undesirable. This problem can be effectively alleviated by leveraging the topological relationship between incidents as done by other approaches. According to Eq. 5, the impact of background noise weakens with distance. However, in FP-Growth, each incident co-occurrence will be counted equally towards the final association rules. UHAS considers the topological similarity by simply calculating the distance. LiDAR employs a more expressive machine learning model, i.e., *node2vec* [22], an algorithmic framework for learning a continuous representation for a network's nodes. However, they both ignore the problem of incomplete failure-impact graph, which is a common issue in online service systems according to our study. The necessity of completing the impact graph will be demonstrated in RQ3. Moreover, different from the traditional applications of graph representation learning, we learn a representation for each unique type of incident, which compactly encodes its relationship with others.

*3) RQ3: The Necessity of the Failure-Impact Graph for Incident Aggregation:* We demonstrate the importance of impact graphs by creating a variant of GRLIA without the phase of failure-impact graph completion (i.e., phase two in Fig. 3), denoted as GRLIA′. We follow LiDAR to remove this feature, which considers two incidents as related only when they are directly connected in the system topology. The experimental results are presented in Table V, where we can see a noticeable drop in the NMI score for all datasets. Due to the high complexity and large scale of online service systems, monitors are often configured in an ad-hoc manner. These monitors may not be able to accommodate to the ever-changing systems and environments. Thus, some incidents are not successfully captured by them. System engineers may

| Method | Dataset1 | Dataset2 | Dataset3 |
|--------|----------|----------|----------|
| GRLIA | **0.831** | **0.866** | **0.912** |
| GRLIA' | 0.782 | 0.808 | 0.846 |

incorrectly perceive the service as healthy, which is a typical situation of gray failures [10]. Without completing the impact graph of failures, the true correlations among incidents cannot be fully recovered.

### D. Threats to Validity

During our study, we have identified the following major threats to the validity.

**Labeling noise**. Our experiments are conducted based on six months of real-world incidents collected from Huawei Cloud. The evaluation requires engineers to inspect and label the incidents manually. Label noises (false positives/false negatives) may be introduced during the manual labeling process. However, the engineers we invite are cloud system professionals and have years of system troubleshooting experience. Moreover, the labeling work can be done quickly and confidently thanks to the incident management system which has user-friendly interfaces. Therefore, we believe the amount of noise is small (if it exists).

**Selection of study subjects**. In our experiments, we only collect incidents from one online service of Huawei Cloud, i.e., the Networking service. This is a large-scale service that supports many upper-layer services such as web application, virtual machine. Sufficient data can be collected from this service system. Another benefit we can enjoy is that the topology of the Networking service system is readily available and accurate. Although we use the Networking service as the subject in this paper, our proposed framework is generalizable, as this service is a typical, representative online service. Thus, we believe GRLIA can be applied to other services and cloud computing platforms and bring them benefits.

The second type of subject that could threaten the validity is the KPI. In production systems, there is a large amount of KPIs available to gauge the similarity between two nodes. Although we only select six representative KPIs (as presented in Section IV-A1), they record the basic and critical states of a service component. Thus, we believe they are able to profile the service system comprehensively.

**Implementation and parameter setting**. The implementation and parameter setting are two critical internal threats to the validity. To reduce the threat of implementation, we employ peer code review. Specifically, the authors are invited to carefully check others' code for mistakes. In terms of parameter setting, we conduct many groups of comparative experiments with different parameters. We choose the parameters by following the original work or empirically based on the best experimental results. In particular, we found GRLIA is not very sensitive to the parameter setting.

## V. Discussion

### A. Success Story

GRLIA has been successfully incorporated into the incident management system of Huawei Cloud. Based on the positive feedback we have received, on-site engineers (OSEs) highly appreciated the novelty of our approach and benefited from it during their daily system maintenance. Specifically, OSEs confirmed the difficulty of the auto-detection of service failures in the existing monitoring system. This is because simple detection techniques (e.g., fixed thresholding) are widely adopted. GRLIA introduces more intelligence and automation by leveraging EVT-based incident burst detection. Interestingly, OSEs found problems for some monitors by comparing their configurations with the aggregated incidents, including wrong names, missing information, etc. Meanwhile, during failure diagnosis, incident aggregation assists OSEs in reducing their investigation scope. Before the deployment of GRLIA, they would have to examine a large number of incidents to locate the failures.

To quantify the practical benefits conveyed to the Networking service system, we further collect failure tickets generated during November 2020. In total, 26 failures are recorded. We calculate the average failure handling time in November and compare it with that in August, September, and October. Results show that the time reduction rate is 24.8%, 21.9%, and 18.6%, respectively, demonstrating the effectiveness of GRLIA in accelerating the incident management of Huawei Cloud.

### B. Lessons Learned

**Optimizing monitor configurations**. Today, popular online services are serving tens of millions of customers. During daily operations, they can produce terabytes and even petabytes of telemetry data such as KPIs, logs, and incidents. However, the majority of these data does not contain much valuable information for service failure analysis. For example, a significant portion of KPIs only record plain system runtime states; most of the incidents are trivial and likely to mitigate automatically with time. The configuration of system monitors should be optimized to report more important yet fewer incidents. In the meantime, monitor configurations show different styles across different service teams, making the monitoring data heterogeneous. Standards should be established for monitor configurations so that high-quality incidents can be created to facilitate the follow-up system analysis, e.g., fault localization.

**Building data collection pipeline**. In online service systems, IT operations play a critical role in system maintenance. Since it is data-driven by nature, modern cloud service providers should build a complete and efficient pipeline for monitoring data collection. Common data quality issues include extremely imbalanced data, small quantity of data, poor signal-to-noise ratio, etc. In general, we are facing the following three challenges: *1) What data should be collected?* We need to identify what metrics and events that are most representative for cloud resource health. Not everything that can be measured needs to be monitored. *2) How to collect*

*and label data?* Labeling incidents (e.g., incident linkages, culprit incidents) requires OSEs to have a decent knowledge about the cloud systems. Since they often devote themselves to emerging issue mitigation and resolution, tools should be developed to facilitate the labeling process, such as label recommendation and friendly interfaces. *3) How to store and query data?* Today's cloud monitoring data are challenging the conventional database systems. To save space, domain-specific compression techniques should be developed, for example, log compression [30]–[32].

## VI. Related Work

### A. Problem Identification

To provide high-quality online services, many researchers have conducted a series of investigations, including problem identification and incident diagnosis from runtime log data and alerts [25], [33], [34]. For example, to identify problems from a large volume of log data, Lin et al. [33] proposed LogCluster to cluster log sequences and pick the center of each cluster. Rosenberg et al. [35] extended LogCluster by incorporating dimension reduction techniques to solve the high-dimension challenge of log sequence vectors. Inspired by LogCluster [33], Zhao et al. [4] clustered online service alerts to identify the representative alerts to engineers. Different from the clustering techniques, Jiang et al. [34] proposed an alert prioritization approach by ranking the importance of alerts based on the KPIs in alert data. The top-ranked alerts are more valuable to identifying problems. However, this approach has a limited scope of application because it is only practical to KPI alerts generated from manually defined threshold rules. To conduct problem identification more aggressively, Chen et al. [36] proposed an incident diagnosis framework to predict general incidents by analyzing their relationships with different alerting signals. Zhao et al. [37] considered a more practical scenario where there are plenty of noisy alerts in online service systems. They proposed eWarn to filter out the noisy alerts and generate interpretable results for incident prediction.

### B. Incident Management

In recent years, cloud computing has gained unprecedented popularity, and incidents are almost inevitable. Thus, incident management becomes a hotspot topic in both academia and industry. Massive amount of effort has been devoted to incident detection [15], [37]–[39] and incident triage [38], [40]–[42]. For example, Lim et al. [43] utilized Hidden Markov Random Field for performance issue clustering to identify representative issues. Chen et al. [42] proposed DeepCT, a deep learning-based approach that is able to accumulate knowledge from incidents' discussions and automate incident triage. However, due to high manual examination costs, these methods cannot handle the overwhelming number of incidents. Many existing work [4], [44] address this problem by reducing the duplicated or correlated alerts. For example, Zhao et al. [45] aimed to recommend the severe alerts to engineers. Lin et al. [39] proposed an alert correlation method to cluster semi-structured alert texts to gain insights from the clustering results.

Similar to our method, Zhao et al. [4] conducted alert reduction by calculating their textual and topological similarity. The centroid alert of each cluster is then selected as the representative incident to engineers. Specifically, they first leveraged conventional methods to detect alert storms and the associated anomalous alerts, and then adopted DBSCAN [46] to cluster alerts based on their textual and topological similarity. Another similar work is LiDAR proposed by Chen et al. [3], which links relevant incidents by incorporating the representation of cloud components. Their framework consists of two modules, a textual encoding module and a component embedding module. The first module learns a representation vector for incident's description in a supervised manner. The textual similarity between two incidents is measured by the cosine distance of their representation vectors. Similarly, the second module learns a vector for system components. The final similarity is calculated by leveraging two parts of information. However, these methods employ a simple weighted sum to combine the information from different sources, and still hardly capture the relationship between incidents. Differently, our method utilizes sophisticated graph representation learning to obtain the semantic relationship of incidents from diverse sources, including temporal locality, topological structure, and KPI metric data. Moreover, many existing incident management methods rely on supervised machine learning techniques to detect anomalies or conduct incident triage. More intelligent approaches with weak-supervision or even unsupervised frameworks are still largely unexplored.

## VII. Conclusion

In this paper, we propose GRLIA, an incident aggregation framework based on graph representation learning. The representation for different types of incidents is learned in an unsupervised and unified fashion, which encodes the interactions among incidents in both temporal and topological dimensions. Online incident aggregation can be efficiently performed by calculating their distance. We have conducted experiments with real-world incidents collected from Huawei Cloud. Compared with fixed thresholding, GRLIA achieves better performance for being able to adjust the threshold automatically. In terms of online incident aggregation, GRLIA also outperforms existing methods by a noticeable margin, confirming its effectiveness. Furthermore, our framework has been successfully incorporated into the incident management system of Huawei Cloud. Feedback from on-site engineers confirms its practical usefulness. We believe our proposed incident aggregation framework can assist engineers in failure understanding and diagnosis.

REFERENCES

[1] S. Wolfe, "Amazon's one hour of downtime on prime day may have cost it up to $100 million in lost sales," 2018. [Online]. Available: https://www.businessinsider.com/amazon-prime-day-website-issues-cost-it-millions-in-lost-sales-2018-7

[2] Z. Chen, Y. Kang, L. Li, X. Zhang, H. Zhang, H. Xu, Y. Zhou, L. Yang, J. Sun, Z. Xu *et al.*, "Towards intelligent incident management: why we need it and how we make it," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 1487–1497.

[3] Y. Chen, X. Yang, H. Dong, X. He, H. Zhang, Q. Lin, J. Chen, P. Zhao, Y. Kang, F. Gao *et al.*, "Identifying linked incidents in large-scale online service systems," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 304–314.

[4] N. Zhao, J. Chen, X. Peng, H. Wang, X. Wu, Y. Zhang, Z. Chen, X. Zheng, X. Nie, G. Wang *et al.*, "Understanding and handling alert storm for online service systems," in *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: Software Engineering in Practice*, 2020, pp. 162–171.

[5] Z. Chen, Y. Kang, F. Gao, L. Yang, J. Sun, Z. Xu, P. Zhao, B. Qiao, L. Li, X. Zhang *et al.*, "Aiops innovations of incident management for cloud services," 2020.

[6] O. Team, "Grlia: Graph-based incident aggregation for large-scale online service systems," 2021. [Online]. Available: https://github.com/OpsPAI/grlia

[7] M. J. Kavis, *Architecting the cloud: design decisions for cloud computing service models (SaaS, PaaS, and IaaS)*. John Wiley & Sons, 2014.

[8] S.-P. Ma, C.-Y. Fan, Y. Chuang, W.-T. Lee, S.-J. Lee, and N.-L. Hsueh, "Using service dependency graph to analyze and test microservices," in *2018 IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC)*, vol. 2. IEEE, 2018, pp. 81–86.

[9] A. Natarajan, P. Ning, Y. Liu, S. Jajodia, and S. E. Hutchinson, *NSDMiner: Automated discovery of network service dependencies*. IEEE, 2012.

[10] P. Huang, C. Guo, L. Zhou, J. R. Lorch, Y. Dang, M. Chintalapati, and R. Yao, "Gray failure: The achilles' heel of cloud-scale systems," in *Proceedings of the 16th Workshop on Hot Topics in Operating Systems*, 2017, pp. 150–155.

[11] W. L. Hamilton, R. Ying, and J. Leskovec, "Representation learning on graphs: Methods and applications," *arXiv preprint arXiv:1709.05584*, 2017.

[12] J. Zhou, G. Cui, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, "Graph neural networks: A review of methods and applications," *arXiv preprint arXiv:1812.08434*, 2018.

[13] A. Siffer, P.-A. Fouque, A. Termier, and C. Largouet, "Anomaly detection in streams with extreme value theory," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017, pp. 1067–1075.

[14] L. De Haan and A. Ferreira, *Extreme value theory: an introduction*. Springer Science & Business Media, 2007.

[15] Q. Lin, K. Hsieh, Y. Dang, H. Zhang, K. Sui, Y. Xu, J.-G. Lou, C. Li, Y. Wu, R. Yao *et al.*, "Predicting node failure in cloud service systems," in *Proceedings of the 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2018, pp. 480–490.

[16] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Soderstrom, "Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding," in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 387–395.

[17] Z. Yang, R. Algesheimer, and C. J. Tessone, "A comparative analysis of community detection algorithms on artificial networks," *Scientific reports*, vol. 6, p. 30750, 2016.

[18] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, p. P10008, 2008.

[19] P. Liu, Y. Chen, X. Nie, J. Zhu, S. Zhang, K. Sui, M. Zhang, and D. Pei, "Fluxrank: A widely-deployable framework to automatically localizing root cause machines for software service failure mitigation," in *2019 IEEE 30th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 2019, pp. 35–46.

[20] E. Keogh and C. A. Ratanamahatana, "Exact indexing of dynamic time warping," *Knowledge and information systems*, vol. 7, no. 3, pp. 358–386, 2005.

[21] J. Han, J. Pei, and Y. Yin, "Mining frequent patterns without candidate generation," *ACM sigmod record*, vol. 29, no. 2, pp. 1–12, 2000.

[22] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, 2016, pp. 855–864.

[23] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2014, pp. 701–710.

[24] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.

[25] S. He, Q. Lin, J.-G. Lou, H. Zhang, M. R. Lyu, and D. Zhang, "Identifying impactful service system problems via log analysis," in *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering - ESEC/FSE 2018*. ACM Press, pp. 60–70.

[26] Stanford, *Evaluation of clustering*, 2008 [Online; accessed November-2020], https://nlp.stanford.edu/IR-book/html/htmledition/evaluation-of-clustering-1.html.

[27] R. Řehůřek, "Gensim: topic modelling for humans," 2009. [Online]. Available: https://radimrehurek.com/gensim/

[28] T. Mikolov, E. Grave, P. Bojanowski, C. Puhrsch, and A. Joulin, "Advances in pre-training distributed word representations," *arXiv preprint arXiv:1712.09405*, 2017.

[29] A. Joulin, E. Grave, P. Bojanowski, M. Douze, H. Jégou, and T. Mikolov, "Fasttext. zip: Compressing text classification models," *arXiv preprint arXiv:1612.03651*, 2016.

[30] J. Liu, J. Zhu, S. He, P. He, Z. Zheng, and M. R. Lyu, "Logzip: extracting hidden structures via iterative clustering for log compression," in *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2019, pp. 863–873.

[31] R. Christensen and F. Li, "Adaptive log compression for massive log data." in *SIGMOD Conference*, 2013, pp. 1283–1284.

[32] P. He, Z. Chen, S. He, and M. R. Lyu, "Characterizing the natural language descriptions in software logging statements," in *2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2018, pp. 178–189.

[33] Q. Lin, H. Zhang, J.-G. Lou, Y. Zhang, and X. Chen, "Log clustering based problem identification for online service systems," in *Proceedings of the 38th International Conference on Software Engineering, ICSE 2016, Austin, TX, USA, May 14-22, 2016 - Companion Volume*. ACM, pp. 102–111.

[34] G. Jiang, H. Chen, K. Yoshihira, and A. Saxena, "Ranking the importance of alerts for problem determination in large computer systems," vol. 14, no. 3, pp. 213–227.

[35] C. M. Rosenberg and L. Moonen, "Improving problem identification via automated log clustering using dimensionality reduction," in *Proceedings of the 12th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. ACM, pp. 1–10.

[36] Y. Chen, X. Yang, Q. Lin, H. Zhang, F. Gao, Z. Xu, Y. Dang, D. Zhang, H. Dong, Y. Xu *et al.*, "Outage prediction and diagnosis for cloud service systems," in *Proceedings of the 2019 International Conference on World Wide Web (WWW)*, 2019, pp. 2659–2665.

[37] N. Zhao, J. Chen, Z. Wang, X. Peng, G. Wang, Y. Wu, F. Zhou, Z. Feng, X. Nie, W. Zhang *et al.*, "Real-time incident prediction for online service systems," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 315–326.

[38] J. Gu, C. Luo, S. Qin, B. Qiao, Q. Lin, H. Zhang, Z. Li, Y. Dang, S. Cai, W. Wu *et al.*, "Efficient incident identification from multi-dimensional issue reports via meta-heuristic search," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 292–303.

[39] D. Lin, R. Raghu, V. Ramamurthy, J. Yu, R. Radhakrishnan, and J. Fernandez, "Unveiling clusters of events for alert and incident management in large-scale enterprise it," in *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New*

*York, NY, USA - August 24 - 27, 2014*, S. A. Macskassy, C. Perlich, J. Leskovec, W. Wang, and R. Ghani, Eds. ACM, pp. 1630–1639.

[40] J. Gao, N. Yaseen, R. MacDavid, F. V. Frujeri, V. Liu, R. Bianchini, R. Aditya, X. Wang, H. Lee, D. Maltz *et al.*, "Scouts: Improving the diagnosis process through domain-customized incident routing," in *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, 2020, pp. 253–269.

[41] J. Chen, X. He, Q. Lin, Y. Xu, H. Zhang, D. Hao, F. Gao, Z. Xu, Y. Dang, and D. Zhang, "An empirical investigation of incident triage for online service systems," in *Proceedings of the 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE Press, 2019, pp. 111–120.

[42] J. Chen, X. He, Q. Lin, H. Zhang, D. Hao, F. Gao, Z. Xu, Y. Dang, and D. Zhang, "Continuous incident triage for large-scale online service systems," in *Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2019, pp. 364–375.

[43] M.-H. Lim, J.-G. Lou, H. Zhang, Q. Fu, A. B. J. Teoh, Q. Lin, R. Ding, and D. Zhang, "Identifying recurrent and unknown performance issues,"

in *2014 IEEE International Conference on Data Mining, ICDM 2014, Shenzhen, China, December 14-17, 2014*, R. Kumar, H. Toivonen, J. Pei, J. Z. Huang, and X. Wu, Eds. IEEE Computer Society, pp. 320–329.

[44] J. Xu, Y. Wang, P. Chen, and P. Wang, "Lightweight and adaptive service api performance monitoring in highly dynamic cloud environment," in *2017 IEEE International Conference on Services Computing, SCC 2017, Honolulu, HI, USA, June 25-30, 2017*, X. F. Liu and U. Bellur, Eds. IEEE Computer Society, pp. 35–43.

[45] N. Zhao, P. Jin, L. Wang, X. Yang, R. Liu, W. Zhang, K. Sui, and D. Pei, "Automatically and adaptively identifying severe alerts for online service systems," in *39th IEEE Conference on Computer Communications, INFOCOM 2020, Toronto, ON, Canada, July 6-9, 2020*. IEEE, pp. 2420–2429.

[46] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), Portland, Oregon, USA*, E. Simoudis, J. Han, and U. M. Fayyad, Eds. AAAI Press, pp. 226–231.