

Debugging Tabular Log as Dynamic Graphs

Anonymous ACL submission

Abstract

Tabular log abstracts objects and events in the real-world system and reports their updates to reflect the change of the system, where one can detect real-world inconsistencies efficiently by debugging corresponding log entries. However, recent advances in processing text-enriched tabular log data overly depend on large language models (LLMs) and other heavy-load models, thus suffering from limited flexibility and scalability. This paper proposes a new framework, GraphLogDebugger, to debug tabular log based on dynamic graphs. By constructing heterogeneous nodes for objects and events and connecting node-wise edges, the framework recovers the system behind the tabular log as an evolving dynamic graph. With the help of our dynamic graph modeling, a simple dynamic Graph Neural Network (GNN) is representative enough to outperform LLMs in debugging tabular log, which is validated by experimental results on real-world log datasets of computer systems and academic papers.

1 Introduction

Tabular log data plays a crucial role in representing and tracking the state and evolution of real-world systems. These logs are structured as rows of log entries, each capturing an event involving certain objects and their attributes at a specific time point. Common examples include system logs recording computing services (Zhu et al., 2023a), research logs tracking scientific publication activities (Clement et al., 2019), and interaction logs from multi-agent systems powered by large language models (LLMs) (Zhang et al., 2025b). Debugging tabular logs is essential: it allows practitioners to detect anomalies in the original systems through efficient inspection of log records.

Log anomaly detection (He et al., 2016) has therefore been a long-standing research field in different niche areas, where data distributions are

invariant or have little change. Existing frameworks (Du et al., 2017; Meng et al., 2019; Zhang et al., 2019; Pei et al., 2020; Guo et al., 2021; Chen and Tsourakakis, 2022) benefit from manually defined data structures or templates for log parsing which are often tailored to certain domains and thus yield absolute success in specific areas like computer system log or financial event log. Due to this domain-specific principle, designing a general-purpose log debugger always remains challenging.

Efforts to overcome this challenge have led to two main lines of work, as shown in Figure 1. One stream focuses on graph modeling of the log data (Cheng et al., 2020; Zehra et al., 2021; Pang et al., 2025), where information in tabular log is gathered in a unified data structure: the graph, such as constructing knowledge graphs or text-rich dynamic graphs for computer system log (Sui et al., 2023; Li et al., 2023). Although these methods are both efficient and powerful, many of them lack flexibility: they still customize static graph structures for certain domains. Another stream explores LLM-based solutions, such as LLM prompting (Yu et al., 2023; Qi et al., 2023; Park, 2024) or retrieval-augmented generation (RAG) (Pan et al., 2024; Zhang et al., 2025a; Wang et al., 2025) pipelines. While these methods demonstrate general capabilities in text-based reasoning, thus showing potential of generalization, they often come with significant drawbacks: high computational costs, slow inference, and difficulty scaling to long log streams or resource-constrained settings.

Inspired by the idea to unify multimodal information in dynamic graphs (Feng et al., 2025), we propose **GraphLogDebugger**, a general and efficient framework for debugging tabular logs through dynamic graph modeling. Our core idea is to interpret tabular log entries as the evolving state of a hidden system, which can be reconstructed as a dynamic heterogeneous graph. We treat objects and events as different types of nodes with text embed-

dings empowered by modern language embedding models, and use the tabular structure to generate time-stamped connections between them. As new log entries arrive, they incrementally update the dynamic graph, capturing both structural and temporal dependencies. This formulation allows us to apply a lightweight dynamic Graph Neural Network (GNN) to perform online anomaly detection by evaluating the likelihood of new connections. Our approach avoids reliance on heavy LLMs while still capturing rich semantic and relational information in the data. Experimental results on real-world datasets from computer system logs and scientific publication logs validate the effectiveness of our approach. Despite its simplicity, our dynamic GNN framework outperforms LLM-based baselines in both accuracy and efficiency, demonstrating that dynamic graph modeling is a highly expressive yet lightweight alternative. Our contributions can be summarized as follows:

- We introduce a novel view of tabular logs as dynamic heterogeneous graphs, bridging the gap between structured attributes and semantic reasoning, and redefine the framework of online log anomaly detection, where object-event connections in each incoming log are evaluated through link prediction.
- We propose a lightweight GNN-based debugger that can efficiently and accurately detect anomalies without using LLMs, and validate its performance on real-world datasets with diverse modalities.

2 Related Works

Tabular Log Processing. Many real-world logs include structured, time-stamped tabular attributes alongside annotated text fields. Examples come from financial prices paired with event series (Tetlock, 2007; Ruiz et al., 2012; Dong et al., 2024), scientific publication metadata (Clement et al., 2019; Kinney et al., 2023), healthcare records (Johnson et al., 2023), computer system logs (Zhu et al., 2023a), and multi-agent system reports (Zhang et al., 2025b). A key challenge in processing tabular logs with machine learning lies in capturing multi-attribute correlations while maintaining comprehension of their semantics (Wu et al., 2025). One common approach integrates main attributes recognized by human priors into structured

data (Yang et al., 2018; Zhao and Feng, 2022; Koval et al., 2024), and then subsequently augments the representation by retrieval (Kurisinkel et al., 2024; Xiao et al., 2025). This modeling achieves good performance in domain-specific data, but lacks flexibility and generalizability for adaptation to other fields (Gardner et al., 2024).

An emerging alternative leverages Large Language Models (LLMs) (Brown et al., 2020), which have demonstrated strong generalizability in understanding, predicting, and generating tabular data (Liu et al., 2023; Zhang et al., 2024b; Fang et al., 2024; Wang et al., 2024b). By parsing diverse logs into a unified format with LLMs (Zhong et al., 2024), these models can be applied to downstream tasks that require reasoning capabilities, such as predicting stock prices (Yu et al., 2023), electricity demand (Wang et al., 2024a), and future events (Shi et al., 2023; Ye et al., 2024). However, LLM-based approaches often suffer from high overheads, complex deployment, and limited throughput. There remains a strong need for lighter-weight alternatives with comparable performance.

Dynamic Graphs. Graph Neural Networks (GNNs) have become a foundational paradigm for learning on graph-structured data (Kipf, 2016; Hamilton et al., 2017; Gilmer et al., 2017). Static GNN models have benefited from advances in message passing (Battaglia et al., 2018), architectural depth (Li et al., 2021; Dwivedi et al., 2020), and inductive scalability (Hamilton et al., 2017). However, many real-world systems are dynamic, motivating models that capture both structural and temporal dependencies. Early approaches used recurrent layers or time-aware embeddings (Li et al., 2017; Seo et al., 2018) to extend static GNNs to dynamic settings (Pareja et al., 2020; Sankar et al., 2020; Kumar et al., 2019). Recent methods have embraced memory modules (Rossi et al., 2020) and temporal encoding (Xu et al., 2020) for finer-grained modeling of time-stamped interactions. Building on this trajectory, ROLAND (You et al., 2022) offers a framework that adapts static GNNs to dynamic graphs via hierarchical state propagation and live-update evaluation, which inspires new advances in benchmarks (Longa et al., 2023; Huang et al., 2023; Zhang et al., 2024a), architectures (Zhu et al., 2023b), explainability (Chen and Ying, 2023), and robustness (Zhang et al., 2023b).

Log Anomaly Detection. Log-based anomaly detection has long been a critical task for sys-

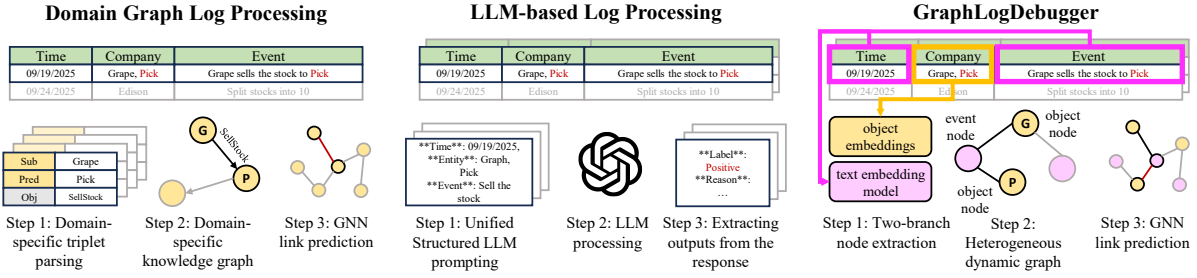


Figure 1: **Comparing GraphLogDebugger with two existing lines of works.** Processing log with domain-specific graphs requires custom text parsing, which lacks flexibility. LLM-based log processing overcomes this shortcoming by the general comprehension skills of LLMs, but suffers from poor efficiency. GraphLogDebugger combines the advantages of graph representation and those of LLMs and balances well generalizability and scalability.

tem reliability, and early neural approaches typically rely on sequence modeling via LSTMs (Du et al., 2017), CNNs (Lu et al., 2018), and autoencoders (Zhang et al., 2021; Castillo et al., 2022; Zhang et al., 2023a). Others incorporate adversarial training (Duan et al., 2021; He et al., 2023), or temporal networks (Zhang et al., 2019; Yang et al., 2021). More recently, pretrained language models have been adopted for log anomaly detection, either via fine-tuning (Guo et al., 2021; Lee et al., 2023) or prompt-based pipelines (Qi et al., 2023; Liu et al., 2024). Retrieval-augmented (No et al., 2024; Pan et al., 2024; Zhang et al., 2025a) methods have further pushed semantic understanding in LLM-based methods. As mentioned, while machine learning-based methods are highly domain-specific, LLM-based methods show some generalizability at a high cost.

One potential solution towards general and scalable methods for log debugging is to introduce dynamic graphs, where tabular log is considered as an evolving system and maintained in a dynamic graph. Early exploration makes use of knowledge graphs (Hogan et al., 2021) with domain specific parsing to generate triplets (Cheng et al., 2020; Zehra et al., 2021; Sui et al., 2023). Recent advances adopt dynamic graphs with text-rich nodes to represent tabular log (Li et al., 2023; Pang et al., 2025). Nevertheless, these works are either domain specific or LLM-based, yet not escaping from the dilemma between generalizability and scalability.

3 Preliminaries

Tabular log is the data modality used to report the update of real-world systems from the perspective of states and relations. It can be formally defined by a time series $X = \{x_0, x_1, x_2, \dots, x_{N-1}\}$ annotated by a timestamps sequence $t_0 < t_1 < t_2 < \dots < t_{N-1}$, where each of x_n is a log en-

try that contains different attributes x_n^m in the table: $x_n = \{x_n^0, x_n^1, x_n^2, \dots, x_n^{M-1}\}$. Summarizing the general case of tabular log data in finance (Dong et al., 2024), healthcare (Johnson et al., 2023), academics (Clement et al., 2019), and other systems (Zhu et al., 2023a), we can separate attributes in the tabular log into two types:

- **Object:** Attributes that represent stand-alone objects in the tabular log, such as companies in the financial news log.
- **Event:** Attributes that describe an event with text, for example, news content in the financial news log and record content in the medical record log. These attributes are usually the center of log entries, where other attributes supplement details and involved objects of the event. Without loss of generality, one log entry only has one Event attribute, because we could merge the text sections of different event attributes into one.

While tabular log abstracts the change of real-world systems, it is expected that we could detect inconsistencies of the system from the corresponding tabular log. Based on the above categorization, we could then define three types of anomalies and corresponding anomaly detection tasks in the tabular log. Give a log entry $x_n = \{x_n^0, x_n^1, x_n^2, \dots, x_n^{M-1}\}$ in the tabular log X :

- **Object anomaly:** Let o_n^p be objects and $\{o_n^0, o_n^1, o_n^2, \dots, o_n^{P-1}\} (P < M)$ be an object set. We have label $y_n = \{y_n^0, y_n^1, y_n^2, \dots, y_n^{P-1}\}$, where $y_n^m = 0$ means that o_n^p is a normal object and $y_n^m = 1$ means that o_n^p is an abnormal object for log entry x_n .
- **Event anomaly:** Given s_n as an event, we have a label y_n , where $y_n = 0$ means that s_n

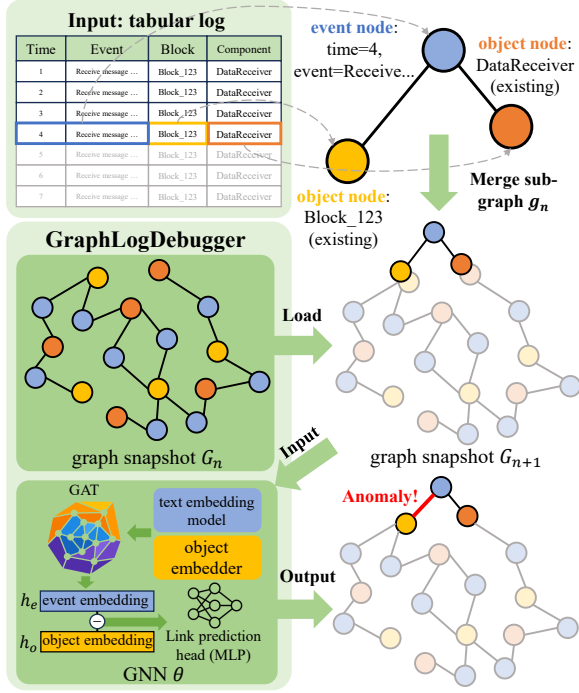


Figure 2: **GraphLogDebugger framework.** The framework checkpoints the GNN θ and the dynamic graph snapshot G_n . When a new log entry emerges, we first extract a sub-graph g_n and use it to update the dynamic graph. Then, we predict the links introduced by g_n , whose results indicate the anomaly.

is a normal event and $y_n = 1$ means that s_n is an abnormal event for log entry x_n .

Considering the tabular log X as an online system where new log entries come dynamically in time order, we could then define the anomaly detection in tabular log as online anomaly detection:

Definition 3.1 (Online Anomaly Detection of Tabular Log). Given an online system that dynamically produces log entry x_n , online anomaly detection for tabular log predicts its anomaly label y_n based on historical log entries $X_n = \{x_0, x_1, x_2, \dots, x_{n-1}\}$

We summarize all used variables in Table 3.

4 GraphLogDebugger

We first integrate tabular log to a heterogeneous dynamic graph (Section 4.1). Then, we reformulate online anomaly detection of tabular log as dynamic graph anomaly detection (Section 4.2). Finally, we apply a dynamic GNN to debug the tabular log (Section 4.3).

4.1 Integrating online tabular log to dynamic graphs

Objects and events in the same log entry are naturally connected in the tabular log, from which we could construct graphs. To this end, we first define the graph structure within one log entry. As shown in Figure 2 (upper section), we build nodes v for both objects and the unique event in the new log entry. Each event and all its objects are connected by an edge e . This yields a sub-graph g_n for each log entry x_n .

Figure 2 (middle section) illustrates the composition of a dynamic graph \mathcal{G} that stores the information of all historical log entries. This dynamic graph gathers all sub-graphs of log entries. We merge identical object nodes so that these sub-graphs are connected. Note that every event node should be unique. Every time a new log entry x_n emerges, we construct a sub-graph g_n accordingly and merge it into the dynamic graph \mathcal{G} . We denote the snapshot of \mathcal{G} at time point t_n by G_n .

4.2 Debugging tabular log graphs as dynamic graphs

Following the above integration, we transfer the online anomaly detection of tabular log defined in Section 3 into an anomaly detection problem of dynamic graph \mathcal{G} (Ekle and Eberle, 2024):

- **Object anomaly detection:** The object anomaly occurs when the edge e between an object node and an event node is abnormal in the latest sub-graph g_n . This anomaly could then be detected by link prediction in the dynamic graph \mathcal{G} (You et al., 2022), where a GNN is applied to predict the likelihood of e . If the likelihood exceeds a threshold, we consider the edge as normal. Otherwise, we consider the edge as an anomaly.
- **Event anomaly detection:** The event anomaly occurs when an abnormal event s is placed in the wrong entry in the latest sub-graph g_n . This means that all edges between this event are anomalies. We can therefore apply link prediction to all edges in the sub-graph and threshold the overall predicted likelihoods to determine the anomaly label.

In a nutshell, the goal of our anomaly detection in Section 3 is equivalent to predicting the likelihood of all edges in the latest log entry sub-graph g_n , based on the dynamic graph snapshot G_n that

stores all historical log entries of tabular log X_n . We finally transfer the online anomaly detection of tabular log into an anomaly detection problem in dynamic graphs, with notations omitted to Table 3:

Definition 4.1 (Online Anomaly Detection of Tabular Log (Dynamic Graph)). Given the snapshot G_n of a dynamic graph $\mathcal{G} = \{G_n\}_{n=0}^{N-1}$ and the new coming sub-graph $g_n = G_{n+1} \setminus G_n$, the goal is to predict the label y_n of links in g_n .

4.3 Designing the GNN for dynamic graph anomaly detection

Figure 2 (bottom section) demonstrates the basic process of using our GNN to predict link anomaly labels. For details, our GNN θ takes the dynamic graph snapshot G_n and the incoming sub-graph g_n as inputs and predicts the likelihood of all links in g_n . The GNN consists of three parts: the node embedder, the GNN backbone, and the prediction head. First, the node embedder offers heterogeneous embeddings for all objects in G_n and g_n and event nodes in G_n . We exclude new events in g_n because we do not expect the outputs of the model will be interfered with by the graph structure in g_n . We assign a unique learnable embedding for each object, and use a pre-trained text embedding model to embed existing events. We also concatenate a time embedding to the event embedding based on the coming time t_n for event s . Our GNN backbone is adapted from the graph attention network (GAT) (Veličković et al., 2017), where we use two separate MLPs to map objects and events to the same space and apply GAT layers for message passing. The prediction head predicts the link between all object-event pairs in g_n . We first compare object embeddings after GAT layers and event text embeddings by reduction. We then pass the result to an MLP with Sigmoid activation to get the likelihood. We omit more details in the design space of the model in Appendix A.3.

Our GNN is trained under the setting of unsupervised anomaly detection (Pang et al., 2021): We separate the dataset into a training split and a test split by chronological order. In the training stage, all links in g_n are normal, and we provide negative examples for training by randomly sampling object-event pairs that are not connected. We append these fraud links to the ground-truth links to balance the label distribution and use them to train the GNN. After backpropagation, we finally update the dynamic graph snapshot G_n with subgraph g_n .

Alg 1 summarizes the training algorithm.

In the test stage, we first resume the GNN as well as the latest dynamic graph snapshot G_n . This achieves the warm start of our debugger system. Then, we construct sub-graphs from the coming log entries and take them as parts of the evolving dynamic graph \mathcal{G} that we succeed from the pre-trained GNN. The evaluation process is summarized in Alg 2.

5 Experiments

5.1 Experimental settings

Datasets. Our work provides a general framework of debugging different types of tabular log under the online setting. To validate this point, we span our experiments over datasets covering four different fields: (1) **Arxiv**: Tabular log recording the timestamps (from 2007-2025), the title and the authors of machine learning papers from the Arxiv (Clement et al., 2019) API; (2) **HDFS**: system log of Hadoop Distributed File System designed to run on commodity hardware (Xu et al., 2009; Zhu et al., 2023a), including the event content together with the objects related to the event; (3) **Analyst**: the commentary records on the finance by analysts, including title, author, and other features of posts ¹; (4) **Landslide**: event catalog reporting the global landslide ². These datasets contain both text-rich attributes and categorical attributes with diversified semantics, thus being challenging to process in one framework efficiently. We omit the details of datasets to Appendix A.1.

We limit the maximum length of all tabular logs to 20,000 by slicing the original datasets. This is because LLM-based baseline methods are costly and not scalable, as discussed in the introduction. To ensure randomness, we randomly pick slices with a length of 20,000 from the whole sliced dataset. For datasets with multiple object attributes, we evaluate object anomaly detection, while for those with only one object attribute, we evaluate event anomaly detection, where event and object anomaly detection are equivalent.

Baselines. Our framework naturally generalizes to tabular log in different domains. Hence, we mainly compare it to baselines which are generally capable of dynamically processing different

¹www.kaggle.com/datasets/miguelaelle/massive-stock-news-analysis-db-for-nlpbacktests

²<https://catalog.data.gov/dataset/global-landslide-catalog-export>

Table 1: **Our proposed GraphLogDebugger outperforms representative baselines on detection effectiveness and efficiency across diverse methods and datasets.** Higher is better for detection effectiveness; lower GFLOPs and higher Throughput are preferred for efficiency. “*” suggests some baselines always predict non-anomaly cases, leading to a 0 prediction, recall, and F1 score.

Method	Dataset: Arxiv Task: Event Anomaly					
	Detection Effectiveness				Efficiency	
	Acc.	Prec.	Recall	F1	GFLOPs	Throughput (it/s)
MLP	0.570 ± 0.205	0.556 ± 0.189	0.893 ± 0.331	0.676 ± 0.028	11.35	825.0 ± 1429.0
RAG (Llama3-70b,k=5)	0.408 ± 0.038	0.426 ± 0.026	0.527 ± 0.029	0.471 ± 0.019	~ 10 ⁵	0.204 ± 0.009
RAG (GPT-oss-20b,k=5)	0.770 ± 0.106	0.771 ± 0.157	0.773 ± 0.014	0.772 ± 0.085	~ 10 ⁴	0.145 ± 0.018
RAG (Llama3-70b,k=10)	0.377 ± 0.014	0.400 ± 0.004	0.493 ± 0.038	0.442 ± 0.014	~ 10 ⁵	0.204 ± 0.015
RAG (GPT-oss-20b,k=10)	0.803 ± 0.090	0.798 ± 0.099	0.813 ± 0.087	0.805 ± 0.087	~ 10 ⁴	0.149 ± 0.014
GraphLogDebugger (Ours)	0.957 ± 0.040	0.920 ± 0.069	1.000 ± 0.000	0.959 ± 0.037	40.39	627.662 ± 7.623
Method	Dataset: Arxiv Task: Object Anomaly					
	Detection Effectiveness				Efficiency	
	Acc.	Prec.	Recall	F1	GFLOPs	Throughput (it/s)
MLP	0.570 ± 0.000	0.538 ± 0.000	0.990 ± 0.000	0.697 ± 0.000	11.35	552.0 ± 167.0
RAG (Llama3-70b,k=5)	0.455 ± 0.033	0.468 ± 0.026	0.667 ± 0.100	0.550 ± 0.052	~ 10 ⁵	0.208 ± 0.018
RAG (GPT-oss-20b,k=5)	0.597 ± 0.019	0.564 ± 0.016	0.850 ± 0.025	0.678 ± 0.004	~ 10 ⁴	0.039 ± 0.018
RAG (Llama3-70b,k=10)	0.463 ± 0.019	0.474 ± 0.011	0.673 ± 0.052	0.556 ± 0.012	~ 10 ⁵	0.154 ± 0.114
RAG (GPT-oss-20b,k=10)	0.598 ± 0.038	0.563 ± 0.023	0.880 ± 0.066	0.687 ± 0.035	~ 10 ⁴	0.039 ± 0.004
GraphLogDebugger (Ours)	0.685 ± 0.065	0.637 ± 0.082	0.870 ± 0.099	0.734 ± 0.024	40.39	592.073 ± 16.513
Method	Dataset: HDFS Task: Object Anomaly					
	Detection Effectiveness				Efficiency	
	Acc.	Prec.	Recall	F1	GFLOPs	Throughput (it/s)
MLP	0.799 ± 0.053	0.801 ± 0.052	0.989 ± 0.000	0.885 ± 0.032	1.4	479.0 ± 194.0
RAG (Llama3-70b,k=5)	0.165 ± 0.022	0 ± 0*	0 ± 0*	0 ± 0*	~ 10 ⁵	0.162 ± 0.085
RAG (GPT-oss-20b,k=5)	0.138 ± 0.029	0 ± 0*	0 ± 0*	0 ± 0*	~ 10 ⁴	0.183 ± 0.010
RAG (Llama3-70b,k=10)	0.173 ± 0.040	0 ± 0*	0 ± 0*	0 ± 0*	~ 10 ⁵	0.194 ± 0.004
RAG (GPT-oss-20b,k=10)	0.138 ± 0.029	0 ± 0*	0 ± 0*	0 ± 0*	~ 10 ⁴	0.192 ± 0.027
GraphLogDebugger (Ours)	0.989 ± 0.023	1.000 ± 0.000	0.987 ± 0.029	0.993 ± 0.015	5.57	529.999 ± 201.308
Method	Dataset: Analyst Task: Event Anomaly					
	Detection Effectiveness				Efficiency	
	Acc.	Prec.	Recall	F1	GFLOPs	Throughput (it/s)
MLP	0.948 ± 0.019	0.922 ± 0.064	0.980 ± 0.043	0.950 ± 0.016	5.58	1996.0 ± 454.0
RAG (Llama3-70b,k=5)	0.408 ± 0.038	0.426 ± 0.026	0.527 ± 0.029	0.471 ± 0.019	~ 10 ⁵	0.204 ± 0.009
RAG (GPT-oss-20b,k=5)	0.770 ± 0.106	0.771 ± 0.157	0.773 ± 0.014	0.772 ± 0.085	~ 10 ⁴	0.145 ± 0.018
RAG (Llama3-70b,k=10)	0.377 ± 0.014	0.400 ± 0.004	0.493 ± 0.038	0.442 ± 0.014	~ 10 ⁵	0.204 ± 0.015
RAG (GPT-oss-20b,k=10)	0.803 ± 0.090	0.798 ± 0.099	0.813 ± 0.087	0.805 ± 0.087	~ 10 ⁴	0.149 ± 0.014
GraphLogDebugger (Ours)	0.957 ± 0.040	0.921 ± 0.069	1.000 ± 0.000	0.959 ± 0.037	8.97	1037.6 ± 286.2
Method	Dataset: Landslide Task: Object Anomaly					
	Detection Effectiveness				Efficiency	
	Acc.	Prec.	Recall	F1	GFLOPs	Throughput (it/s)
MLP	0.831 ± 0.079	0.842 ± 0.180	0.841 ± 0.149	0.838 ± 0.056	5.58	5391.0 ± 328.0
RAG (Llama3-70b,k=5)	0.543 ± 0.074	0.944 ± 0.056	0.095 ± 0.156	0.168 ± 0.267	~ 10 ⁵	0.355 ± 0.044
RAG (GPT-oss-20b,k=5)	0.611 ± 0.068	0.701 ± 0.050	0.389 ± 0.246	0.495 ± 0.195	~ 10 ⁴	0.155 ± 0.136
RAG (Llama3-70b,k=10)	0.551 ± 0.034	0.904 ± 0.204	0.119 ± 0.102	0.208 ± 0.160	~ 10 ⁵	0.321 ± 0.193
RAG (GPT-oss-20b,k=10)	0.623 ± 0.074	0.723 ± 0.081	0.397 ± 0.149	0.511 ± 0.144	~ 10 ⁴	0.166 ± 0.008
GraphLogDebugger (Ours)	0.840 ± 0.080	0.798 ± 0.117	0.929 ± 0.059	0.858 ± 0.062	19.70	1334.13 ± 108.40

types of tabular log that contains text-rich and categorical attributes. **MLP** exploits a pretrained text embedding model to embed events and a learnable embedding for objects. A 3-layer MLP is then applied to map the embeddings to anomaly scores. We also compare a series of baselines based on retrieval augmented generation (**RAG**) (Lewis et al., 2020), which is the mainstream method to process general tabular log in the realistic scenario (Akhtar et al., 2025). We deploy RAG based on two advanced open-sourced LLMs, Llama-3-70b (Dubey et al., 2024) and GPT-oss-20b (Agarwal et al., 2025) with the 5 and 10 retrieval entries. The retrieval database is built on the whole training split and the seen log entries during the online evaluation. We construct specified prompts for different datasets and omitted the description to Appendix A. Both MLP and all RAG baselines use all-MiniLM-L6-v2³ (Reimers and Gurevych, 2019) as the text embedding model. We do not compare to baselines on log anomaly detection because all these methods are either template-based or domain-specific, which cannot be applied to datasets other than computer system log.

Task. Following the setting of unsupervised anomaly detection (Liu et al., 2021; Schmidl et al., 2022), our basic task is to output an anomaly score for each log entry x_n at timestamp t_n , where higher scores denote more outlyingness (Han et al., 2022). In our task, we use 1 to denote anomalies and 0 to denote normal examples in the ground-truth. Following the setup of online anomaly detection (Qiao et al., 2025), we use the first 90% split of the dataset for training, where both the log entries and their anomaly labels are available to access for methods. Methods train the model on this training split or use it for the retrieval database. For the rest 10%, we use it as the test split in our online evaluation, where methods can make use of the seen log entries but their anomaly labels are not accessible. We study two types of anomalies in our experiments: object anomalies and event anomalies. Following the definition in Section 3, we inject object anomalies by swapping an object in the log entry with another existing object. To ensure that historical data contains useful information, we only perturb existing objects in the history. Event anomalies are generated similarly by swapping events. The anomaly rate is set to be 0.05.

³<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

Evaluation. We calculate the metrics for information retrieval: accuracy, precision, recall, and f1 score for the dataset. We also evaluate the efficiency of different methods by GFlops and throughputs. During evaluation, we notice that LLM-based baselines tend to be very slow in processing speed. Hence, we include all anomaly log entries and 50 random normal entries in a subset and run RAG only on this subset. For other baselines and our methods, we obtain the prediction result for the full test split but only compute the metric on the above subset for fair comparison. We run experiments three times and post the average value of metrics with error bars.

GraphLogDebugger. The GNN architecture in GraphLogDebugger is a 3-layer GAT backbone with a two-branch node embedder and an MLP prediction head. The node embedder uses 512-d embeddings for objects and the text embedding of all-MiniLM-L6-v2 (Reimers and Gurevych, 2019) for events, with an MLP to map them into the same space. The embedding size of GAT and the prediction head is also 512. We train the GNN for 10 epochs under the learning rate 0.0001 on Adam and the negative ratio 10 on the training split. Following You et al. (2022), we set the batch size as 1 and use a window length of 100.

5.2 Main Results

Table 1 shows that **GraphLogDebugger** consistently outperforms both MLP and RAG-based baselines across all tabular log datasets on five tasks, in terms of detection performance and efficiency.

Effectiveness: GraphLogDebugger achieves the highest F1 scores across all tasks, outperforming RAG baselines—especially in structurally complex domains like HDFS, where RAG methods fail to detect meaningful anomalies (F1 = 0.0). Specifically, RAG baselines are completely fooled by the anomaly pattern that their predicted labels depend on whether there is an existing record with the same format in the retrieved examples, which does not contribute to a reasonable prediction. Two tasks on the Arxiv dataset are the most difficult, where GraphLogDebugger still beats baselines with a higher precision in not abusing anomaly prediction. Even in semantically rich settings such as Analyst and LandSlide, where RAG baselines are expected to excel, our model surpasses them.

Efficiency: RAG approaches exhibit extremely low throughput (typically below 0.3 iterations per second) due to the computational overhead of large

language models. In contrast, GraphLogDebugger achieves throughput of at least 500 per second, with significantly lower GFLOPs, enabling real-time anomaly detection in high-throughput environments.

5.3 Case study: Where does RAG fail?

It is natural that GraphLogDebugger yields advantages in efficiency compared to the RAG-baseline, for the latter relies on LLMs with billions of parameters. However, the leading performance of GraphLogDebugger in detection needs further explanation, while RAG enjoys the general comprehension and reasoning ability of modern LLMs. To this end, we study cases from event anomaly detection of the Arxiv dataset. We choose this task because the degree of event nodes can directly reflect the local graph density of the node-of-interest. We calculate the correlation between event node degrees and the accuracy of GraphLogDebugger and that of RAG(GPT-oss-20b, $k=10$). Quantitatively, the correlation between the accuracy of GraphLogDebugger and the node degree is 0.1561, while the accuracy of RAG is negatively correlated with the node degree with the correlation of -0.1087. This indicates that GraphLogDebugger outperforms RAG on event nodes with rich connections with objects, where semantics of these objects are necessary to detect the anomaly.

We further raise three cases to investigate when and how GraphLogDebugger and RAG fail. In Case 1, RAG predicts the normal example as abnormal because the limited retrieved examples do not provide enough evidence to prove the coherence of the author team. By contrast, GraphLogDebugger validates overall team consistency by checking the research background of every author, which correctly predicts the negative label. Case 2 is complementary to Case 1, where GraphLogDebugger is able to scan the research interest of every author and detect the anomaly accurately. However, when the connected objects are few, such as in Case 3, GraphLogDebugger may not have enough references based on the graph to make a correct judgment. In similar cases, RAG could then outperform GraphLogDebugger to recognize patterns in the number of authors in the same domain.

These cases provide insights on how graphs can benefit retrieval augmented generation. When the key entry has dense connections with other entries, traditional retrieval based on similarity cannot efficiently include enough entries to enhance the gener-

ation quality. With the help of modern embedding models, graphs can be introduced to gather information in these multi-entry scenarios.

Case 1: Label=**negative**, RAG=**positive**, GraphLogDebugger=**negative**

Title: "SymbioSim: Human-in-the-loop Simulation Platform for Bidirectional Continuing Learning in Human-Robot Interaction"

Authors: "C1_A1", "C1_A2", "C1_A3", "C1_A4", "C1_A5", "C1_A6", "C1_A7", "C1_A8", "C1_A9" (9 objects)

Reason (RAG): The author team composition, research domain mismatch, and unclear collaboration patterns raise suspicions about the coherence of the record.

Case 2: Label=**positive**, RAG=**negative**, GraphLogDebugger=**positive**

Title: "VERA: Explainable Video Anomaly Detection via Verbalized Learning of Vision-Language Models"

Authors: "C2_A1", "C2_A2", "C2_A3", "C2_A4", "C2_A5", "C2_A6", "C2_A7" (7 objects)

Reason (RAG): The record seems coherent, with individual authors' expertise areas aligning with the paper's topic, although the team size is slightly larger than expected.

Case 3: Label=**positive**, RAG=**positive**, GraphLogDebugger=**negative**

Title: "Transformer⁻¹: Input-Adaptive Computation for Resource-Constrained Deployment"

Authors: "C3_A1" (1 objects)

Reason (RAG): The record consists of a single author, which is consistent with similar papers in the same research domain.

6 Conclusion

We propose a general framework to cover online debugging for heterogeneous tabular logs. By modeling online log debugging as anomaly detection of dynamic graphs, our framework integrates different types of log data into a unified modality by text embedding models, where a dynamic GNN debugs the log through link prediction. Our framework shows good performance in four different datasets while maintaining high efficiency compared to the mainstream RAG-based method. In the future, we will focus on bug correction based on the current bug detection system.

588 Limitations

589 Our work explores combining dynamic GNNs and
590 text embedding models to process log data under
591 the online setting, which indicates the potential to
592 accelerate the online process of data streams by
593 a graph-based method. Nevertheless, our experi-
594 ments mainly show this potential in the bug detec-
595 tion setting. We leave the exploration of online bug
596 correction to future work.

597 Ethical considerations

598 Our work focuses on detecting inconsistencies
599 in general tabular log data, which enhances the
600 progress of automated log data processing in real-
601 world scenarios. While automation of log process-
602 ing may raise issues concerning hallucination or
603 fraud reporting, our work does not explicitly intro-
604 duce new risks compared to existing research.

605 References

606 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Alt-
607 man, Andy Applebaum, Edwin Arbus, Rahul K
608 Arora, Yu Bai, Bowen Baker, Haiming Bao, and 1
609 others. 2025. gpt-oss-120b & gpt-oss-20b model
610 card. *arXiv preprint arXiv:2508.10925*.

611 Siraj Akhtar, Saad Khan, and Simon Parkinson. 2025.
612 Llm-based event log analysis techniques: A survey.
613 *arXiv preprint arXiv:2502.00677*.

614 Peter W Battaglia, Jessica B Hamrick, Victor Bapst,
615 Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Ma-
616 teusz Malinowski, Andrea Tacchetti, David Raposo,
617 Adam Santoro, Ryan Faulkner, and 1 others. 2018.
618 Relational inductive biases, deep learning, and graph
619 networks. *arXiv preprint arXiv:1806.01261*.

620 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
621 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
622 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
623 Askeel, and 1 others. 2020. Language models are
624 few-shot learners. *Advances in neural information
625 processing systems*, 33:1877–1901.

626 Martina Castillo, Antonio Pecchia, and Ugo Villano.
627 2022. Autolog: Anomaly detection by deep autoen-
628 coding of system logs. *Expert Systems with Applica-
629 tions*, 191.

630 Jialin Chen and Rex Ying. 2023. Tempme: Towards
631 the explainability of temporal graph neural networks
632 via motif discovery. *Advances in Neural Information
633 Processing Systems*, 36:29005–29028.

634 Tianyi Chen and Charalampos Tsourakakis. 2022. An-
635 tibenford subgraphs: Unsupervised anomaly detec-
636 tion in financial networks. In *Proceedings of the 28th
637 ACM SIGKDD Conference on Knowledge Discovery
638 and Data Mining*, pages 2762–2770.

Dawei Cheng, Fangzhou Yang, Xiaoyang Wang, Ying
Zhang, and Liqing Zhang. 2020. Knowledge graph-
based event embedding framework for financial quan-
titative investments. In *Proceedings of the 43rd In-
ternational ACM SIGIR Conference on Research and
Development in Information Retrieval*, pages 2221–
2230.

Colin B Clement, Matthew Bierbaum, Kevin P
O’Keeffe, and Alexander A Alemi. 2019. On
the use of arxiv as a dataset. *arXiv preprint
arXiv:1905.00075*.

Zihan Dong, Xinyu Fan, and Zhiyuan Peng. 2024. **Fn-
spid: A comprehensive financial news dataset in time
series**. *Preprint*, arXiv:2402.06698.

Min Du, Feifei Li, Guineng Zheng, and Vivek Srikumar.
2017. Deeplog: Anomaly detection and diagnosis
from system logs through deep learning. In *Pro-
ceedings of the 2017 ACM SIGSAC conference on
computer and communications security*, pages 1285–
1298.

Xiaoyu Duan, Shi Ying, Wanli Yuan, Hailong Cheng,
and Xiang Yin. 2021. A generative adversarial net-
works for log anomaly detection. *Computer Systems
Science & Engineering*, 37(1).

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, and 1 others. 2024. The llama 3 herd of models.
arXiv e-prints, pages arXiv–2407.

Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas
Laurent, Yoshua Bengio, and Xavier Bresson. 2020.
Benchmarking graph neural networks. *arXiv preprint
arXiv:2003.00982*.

Ocheme Anthony Ekle and William Eberle. 2024.
Anomaly detection in dynamic graphs: A compre-
hensive survey. *ACM Transactions on Knowledge
Discovery from Data*, 18(8):1–44.

Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang,
Ziqing Hu, Yanjun Qi, Scott Nickleach, Diego Socol-
insky, Srinivasan Sengamedu, and Christos Faloutsos.
2024. Large language models (llms) on tabular data:
Prediction, generation, and understanding—a survey.
arXiv preprint arXiv:2402.17944.

Tao Feng, Yexin Wu, Guanyu Lin, and Jiaxuan
You. 2025. Graph world model. *arXiv preprint
arXiv:2507.10539*.

Josh Gardner, Juan C Perdomo, and Ludwig Schmidt.
2024. Large scale transfer learning for tabular data
via language modeling. *Advances in Neural Informa-
tion Processing Systems*, 37:45155–45205.

Justin Gilmer, Samuel S Schoenholz, Patrick F Riley,
Oriol Vinyals, and George E Dahl. 2017. Neural mes-
sage passing for quantum chemistry. In *International
conference on machine learning*, pages 1263–1272.
Pmlr.

694	Haixuan Guo, Shuhan Yuan, and Xintao Wu. 2021.	<i>ACM SIGKDD international conference on knowl-</i>	749
695	Logbert: Log anomaly detection via bert. In <i>2021</i>	<i>edge discovery & data mining</i> , pages 1269–1278.	750
696	<i>international joint conference on neural networks</i>		
697	(<i>IJCNN</i>), pages 1–8. IEEE.		
698	Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017.	Litton Jose Kurisinkel, Pruthwik Mishra, and Yue	751
699	Inductive representation learning on large graphs. <i>Ad-</i>	Zhang. 2024. Text2timeseries: Enhancing financial	752
700	<i>vances in neural information processing systems</i> , 30.	forecasting through time series prediction updates	753
		with event-driven insights from large language mod-	754
		els. <i>arXiv preprint arXiv:2407.03689</i> .	755
701	Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi	Yukyung Lee, Jina Kim, and Pilsung Kang. 2023.	756
702	Jiang, and Yue Zhao. 2022. Adbench: Anomaly de-	Lanobert: System log anomaly detection based on	757
703	tection benchmark. <i>Advances in neural information</i>	bert masked language model. <i>Applied Soft Comput-</i>	758
704	<i>processing systems</i> , 35:32142–32159.	<i>ing</i> , 146:110689.	759
705	Shilin He, Jieming Zhu, Pinjia He, and Michael R Lyu.	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	760
706	2016. Experience report: System log analysis for	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	761
707	anomaly detection. In <i>2016 IEEE 27th international</i>	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	762
708	<i>symposium on software reliability engineering (IS-</i>	täschel, and 1 others. 2020. Retrieval-augmented gen-	763
709	<i>SRE</i>), pages 207–218. IEEE.	eration for knowledge-intensive nlp tasks. <i>Advances</i>	764
		<i>in neural information processing systems</i> , 33:9459–	765
710	Zhangyue He, Yanni Tang, Kaiqi Zhao, Jiamou Liu, and	9474.	766
711	Wu Chen. 2023. Graph-based log anomaly detec-		
712	tion via adversarial training. In <i>International Sympo-</i>	Guohao Li, Matthias Müller, Ali Thabet, and Bernard	767
713	<i>sium on Dependable Software Engineering: Theories,</i>	Ghanem. 2021. Deepgcns: Can gcns go as deep as	768
714	<i>Tools, and Applications</i> , pages 55–71. Springer.	cnns? <i>IEEE Transactions on Pattern Analysis and</i>	769
		<i>Machine Intelligence</i> .	770
715	Aidan Hogan, Eva Blomqvist, Michael Cochez, Clau-	Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu.	771
716	dia d’Amato, Gerard De Melo, Claudio Gutierrez,	2017. Diffusion convolutional recurrent neural net-	772
717	Sabrina Kirrane, José Emilio Labra Gayo, Roberto	work: Data-driven traffic forecasting. <i>arXiv preprint</i>	773
718	Navigli, Sebastian Neumaier, and 1 others. 2021.	<i>arXiv:1707.01926</i> .	774
719	Knowledge graphs. <i>ACM Computing Surveys (Csur)</i> ,		
720	54(4):1–37.	Yufei Li, Yanchi Liu, Haoyu Wang, Zhengzhang Chen,	775
		Wei Cheng, Yuncong Chen, Wenchao Yu, Haifeng	776
721	Shenyang Huang, Farimah Poursafaei, Jacob Danovitch,	Chen, and Cong Liu. 2023. Glad: Content-aware	777
722	Matthias Fey, Weihua Hu, Emanuele Rossi, Jure	dynamic graphs for log anomaly detection. In <i>2023</i>	778
723	Leskovec, Michael Bronstein, Guillaume Rabusseau,	<i>IEEE International Conference on Knowledge Graph</i>	779
724	and Reihaneh Rabbany. 2023. Temporal graph bench-	(<i>ICKG</i>), pages 9–18. IEEE.	780
725	mark for machine learning on temporal graphs. <i>Ad-</i>		
726	<i>vances in Neural Information Processing Systems</i> ,	Tianyang Liu, Fei Wang, and Muhao Chen. 2023. Re-	781
727	36:2056–2073.	thinking tabular data understanding with large lan-	782
		guage models. <i>arXiv preprint arXiv:2312.16702</i> .	783
728	Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin	Yilun Liu, Shimin Tao, Weibin Meng, Feiyu Yao, Xi-	784
729	Gayles, Ayad Shammout, Steven Horng, Tom J Pol-	aofeng Zhao, and Hao Yang. 2024. Logprompt:	785
730	lard, Sicheng Hao, Benjamin Moody, Brian Gow, and	Prompt engineering towards zero-shot and inter-	786
731	1 others. 2023. Mimic-iv, a freely accessible elec-	pretable log analysis. In <i>Proceedings of the 2024</i>	787
732	tronic health record dataset. <i>Scientific data</i> , 10(1):1.	<i>IEEE/ACM 46th International Conference on Soft-</i>	788
		<i>ware Engineering: Companion Proceedings</i> , pages	789
733	Rodney Kinney, Chloe Anastasiades, Russell Authur,	364–365.	790
734	Iz Beltagy, Jonathan Bragg, Alexandra Buraczyn-	Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan	791
735	ski, Isabel Cachola, Stefan Candra, Yoganand Chan-	Zhou, and George Karypis. 2021. Anomaly detection	792
736	drasekhar, Arman Cohan, and 1 others. 2023. The	on attributed networks via contrastive self-supervised	793
737	semantic scholar open data platform. <i>arXiv preprint</i>	learning. <i>IEEE transactions on neural networks and</i>	794
738	<i>arXiv:2301.10140</i> .	<i>learning systems</i> , 33(6):2378–2392.	795
739	TN Kipf. 2016. Semi-supervised classification with	Alessandro Longa, Valerio Lachi, Giovanni Santin,	796
740	graph convolutional networks. <i>arXiv preprint</i>	Monica Bianchini, Bruno Lepri, and 1 others. 2023.	797
741	<i>arXiv:1609.02907</i> .	Graph neural networks for temporal graphs: State	798
		of the art, open challenges, and opportunities. <i>arXiv</i>	799
742	Ross Koval, Nicholas Andrews, and Xifeng Yan. 2024.	<i>preprint arXiv:2305.12472</i> .	800
743	Financial forecasting from textual and tabular time		
744	series. In <i>Findings of the Association for Computa-</i>	Siyang Lu, Xiang Wei, Yandong Li, and Liqiang Wang.	801
745	<i>tional Linguistics: EMNLP 2024</i> , pages 8289–8300.	2018. Detecting anomaly in big data system logs	802
		using convolutional neural network. In <i>2018 IEEE</i>	803
746	Srijan Kumar, Xikun Zhang, and Jure Leskovec. 2019.		
747	Predicting dynamic embedding trajectory in tempo-		
748	ral interaction networks. In <i>Proceedings of the 25th</i>		

804		16th Intl Conf on Dependable, Autonomic and Secure Computing, 16th Intl Conf on Pervasive Intelligence and Computing, 4th Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech), pages 151–158. IEEE.	
810	Weibin Meng, Ying Liu, Yichen Zhu, Shenglin Zhang, Dan Pei, Yuqing Liu, Yihao Chen, Ruizhi Zhang, Shimin Tao, Pei Sun, and 1 others. 2019. Loganomaly: Unsupervised detection of sequential and quantitative anomalies in unstructured logs. In <i>IJCAI</i> , volume 19, pages 4739–4745.		
816	Gunho No, Yukyung Lee, Hyeongwon Kang, and Pilsung Kang. 2024. Training-free retrieval-based log anomaly detection with pre-trained language model considering token-level information. <i>Engineering Applications of Artificial Intelligence</i> , 133:108613.		
821	Jonathan Pan, Wong Swee Liang, and Yuan Yidi. 2024. Raglog: Log anomaly detection using retrieval augmented generation. In <i>2024 IEEE World Forum on Public Safety Technology (WFPST)</i> , pages 169–174. IEEE.		
826	Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. 2021. Deep learning for anomaly detection: A review. <i>ACM computing surveys (CSUR)</i> , 54(2):1–38.		
830	Yunhe Pang, Bo Chen, Fanjin Zhang, Yanghui Rao, Evgeny Kharlamov, and Jie Tang. 2025. Guard: Effective anomaly detection through a text-rich and graph-informed language model. In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2</i> , pages 2222–2233.		
836	Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao Schardl, and Charles Leiserson. 2020. Evolvegc: Evolving graph convolutional networks for dynamic graphs. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 5363–5370.		
843	Taejin Park. 2024. Enhancing anomaly detection in financial markets with an llm-based multi-agent framework. <i>arXiv preprint arXiv:2403.19735</i> .		
846	Yulong Pei, Fang Lyu, Werner Van Ipenburg, and Mykola Pechenizkiy. 2020. Subgraph anomaly detection in financial transaction networks. In <i>Proceedings of the First ACM International Conference on AI in Finance</i> , pages 1–8.		
851	Jiaxing Qi, Shaohan Huang, Zhongzhi Luan, Shu Yang, Carol Fung, Hailong Yang, Depei Qian, Jing Shang, Zhiwen Xiao, and Zhihui Wu. 2023. Loggpt: Exploring chatgpt for log-based anomaly detection. In <i>2023 IEEE International Conference on High Performance Computing & Communications, Data Science & Systems, Smart City & Dependability in Sensor, Cloud & Big Data Systems & Application (HPCC/DSS/SmartCity/DependSys)</i> , pages 273–280. IEEE.		
	Hezhe Qiao, Hanghang Tong, Bo An, Irwin King, Charu Aggarwal, and Guansong Pang. 2025. Deep graph anomaly detection: A survey and new perspectives. <i>IEEE Transactions on Knowledge and Data Engineering</i> .		860 861 862 863 864
	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.		865 866 867 868 869
	Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. 2020. Temporal graph networks for deep learning on dynamic graphs. <i>arXiv preprint arXiv:2006.10637</i> .		870 871 872 873
	Eduardo J Ruiz, Vagelis Hristidis, Carlos Castillo, Aristides Gionis, and Alejandro Jaimes. 2012. Correlating financial time series with micro-blogging activity. In <i>Proceedings of the fifth ACM international conference on Web search and data mining</i> , pages 513–522.		874 875 876 877 878
	Aravind Sankar, Yanhong Wu, Liang Gou, Wei Zhang, and Hao Yang. 2020. Dysat: Deep neural representation learning on dynamic graphs via self-attention networks. In <i>Proceedings of the 13th international conference on web search and data mining</i> , pages 519–527.		879 880 881 882 883 884
	Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. 2022. Anomaly detection in time series: a comprehensive evaluation. <i>Proceedings of the VLDB Endowment</i> , 15(9):1779–1797.		885 886 887 888
	Youngjoo Seo, Michaël Defferrard, Pierre Vandergheynst, and Xavier Bresson. 2018. Structured sequence modeling with graph convolutional recurrent networks. In <i>International conference on neural information processing</i> , pages 362–373. Springer.		889 890 891 892 893
	Xiaoming Shi, Siqiao Xue, Kangrui Wang, Fan Zhou, James Zhang, Jun Zhou, Chenhao Tan, and Hongyuan Mei. 2023. Language models can improve event prediction by few-shot abductive reasoning. <i>Advances in Neural Information Processing Systems</i> , 36:29532–29557.		894 895 896 897 898 899
	Yicheng Sui, Yuzhe Zhang, Jianjun Sun, Ting Xu, Shenglin Zhang, Zhengdan Li, Yongqian Sun, Fangrui Guo, Junyu Shen, Yuzhi Zhang, and 1 others. 2023. Logkg: Log failure diagnosis through knowledge graph. <i>IEEE Transactions on Services Computing</i> , 16(5):3493–3507.		900 901 902 903 904 905
	Paul C Tetlock. 2007. Giving content to investor sentiment: The role of media in the stock market. <i>The Journal of finance</i> , 62(3):1139–1168.		906 907 908
	Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. <i>arXiv preprint arXiv:1710.10903</i> .		909 910 911 912

913	Jingru Wang, Wen Ding, and Xiaotong Zhu. 2025.	Xinli Yu, Zheng Chen, Yuan Ling, Shujing Dong,	967
914	Financial analysis: Intelligent financial data anal-	Zongyi Liu, and Yanbin Lu. 2023. Temporal data	968
915	ysis system based on llm-rag. <i>arXiv preprint</i>	meets llm-explainable financial time series forecast-	969
916	<i>arXiv:2504.06279</i> .	ing. <i>arXiv preprint arXiv:2306.11025</i> .	970
917	Xinlei Wang, Maike Feng, Jing Qiu, Jinjin Gu, and Jun-	Samreen Zehra, Syed Farhan Mohsin Mohsin, Shaukat	971
918	hua Zhao. 2024a. From news to forecast: Integrating	Wasi, Syed Imran Jami, Muhammad Shoaib Sid-	972
919	event analysis in llm-based time series forecasting	diqui, and Muhammad Khaliq-Ur-Rahman Raazi	973
920	with reflection. <i>Advances in Neural Information Pro-</i>	Syed. 2021. Financial knowledge graph based fi-	974
921	<i>cessing Systems</i> , 37:58118–58153.	nancial report query system. <i>IEEE Access</i> , 9:69766–	975
922	Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Mar-	69782.	976
923	tin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly	Jiasheng Zhang, Jialin Chen, Menglin Yang, Aosong	977
924	Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu	Feng, Shuang Liang, Jie Shao, and Rex Ying. 2024a.	978
925	Lee, and 1 others. 2024b. Chain-of-table: Evolving	Dtgb: A comprehensive benchmark for dynamic text-	979
926	tables in the reasoning chain for table understanding.	attributed graphs. <i>Advances in Neural Information</i>	980
927	<i>arXiv preprint arXiv:2401.04398</i> .	<i>Processing Systems</i> , 37:91405–91429.	981
928	Xiaofeng Wu, Alan Ritter, and Wei Xu. 2025. Tab-	Lingzhe Zhang, Tong Jia, Mengxi Jia, Yifan Wu,	982
929	ular data understanding with llms: A survey of	Hongyi Liu, and Ying Li. 2025a. Xraglog:	983
930	recent advances and challenges. <i>arXiv preprint</i>	A resource-efficient and context-aware log-based	984
931	<i>arXiv:2508.00217</i> .	anomaly detection method using retrieval-augmented	985
932	Mengxi Xiao, Zihao Jiang, Lingfei Qian, Zhengyu Chen,	generation. In <i>AAAI 2025 Workshop on Preventing</i>	986
933	Yueru He, Yijing Xu, Yuecheng Jiang, Dong Li,	<i>and Detecting LLM Misinformation (PDLM)</i> .	987
934	Ruey-Ling Weng, Min Peng, and 1 others. 2025.	Linming Zhang, Wenzhong Li, Zhijie Zhang, Qingning	988
935	Enhancing financial time-series forecasting with	Lu, Ce Hou, Peng Hu, Tong Gui, and Sanglu Lu.	989
936	retrieval-augmented large language models. <i>arXiv</i>	2021. Logattn: Unsupervised log anomaly detection	990
937	<i>preprint arXiv:2503.67890</i> .	with an autoencoder based attention mechanism. In	991
938	Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant	<i>International conference on knowledge science, engi-</i>	992
939	Kumar, and Kannan Achan. 2020. Inductive repre-	<i>neering and management</i> , pages 222–235. Springer.	993
940	sentation learning on temporal graphs. <i>arXiv preprint</i>	Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu,	994
941	<i>arXiv:2002.07962</i> .	Zhiguang Han, Jingyang Zhang, Beibin Li, Chi	995
942	Wei Xu, Ling Huang, Armando Fox, David Patterson,	Wang, Huazheng Wang, Yiran Chen, and 1 others.	996
943	and Michael I Jordan. 2009. Detecting large-scale	2025b. Which agent causes task failures and when?	997
944	system problems by mining console logs. In <i>Pro-</i>	on automated failure attribution of llm multi-agent	998
945	<i>ceedings of the ACM SIGOPS 22nd symposium on</i>	systems. <i>arXiv preprint arXiv:2505.00212</i> .	999
946	<i>Operating systems principles</i> , pages 117–132.	Xiaokang Zhang, Sijia Luo, Bohan Zhang, Zeyao	1000
947	Hang Yang, Yubo Chen, Kang Liu, Yang Xiao, and Jun	Ma, Jing Zhang, Yang Li, Guanlin Li, Zijun Yao,	1001
948	Zhao. 2018. Dcfee: A document-level chinese finan-	Kangli Xu, Jinchang Zhou, and 1 others. 2024b.	1002
949	cial event extraction system based on automatically	Tablellm: Enabling tabular data manipulation by	1003
950	labeled training data. In <i>Proceedings of ACL 2018,</i>	llms in real office usage scenarios. <i>arXiv preprint</i>	1004
951	<i>System Demonstrations</i> , pages 50–55.	<i>arXiv:2403.19318</i> .	1005
952	Lin Yang, Junjie Chen, Zan Wang, Weijing Wang, Jia-	Xinye Zhang, Xiaoli Chai, Minghua Yu, and Ding Qiu.	1006
953	jun Jiang, Xuyuan Dong, and Wenbin Zhang. 2021.	2023a. Anomaly detection model for log based on	1007
954	Semi-supervised log-based anomaly detection via	lstm network and variational autoencoder. In <i>2023</i>	1008
955	probabilistic label estimation. In <i>2021 IEEE/ACM</i>	<i>4th International Conference on Information Science,</i>	1009
956	<i>43rd International Conference on Software Engineer-</i>	<i>Parallel and Distributed Systems (ISPDS)</i> , pages 239–	1010
957	<i>ing (ICSE)</i> , pages 1448–1460. IEEE.	244. IEEE.	1011
958	Chenchen Ye, Ziniu Hu, Yihe Deng, Zijie Huang,	Xu Zhang, Yong Xu, Qingwei Lin, Bo Qiao, Hongyu	1012
959	Mingyu Derek Ma, Yanqiao Zhu, and Wei Wang.	Zhang, Yingnong Dang, Chunyu Xie, Xinsheng Yang,	1013
960	2024. Mirai: Evaluating llm agents for event fore-	Qian Cheng, Ze Li, and 1 others. 2019. Robust log-	1014
961	casting. <i>arXiv preprint arXiv:2407.01231</i> .	based anomaly detection on unstable log data. In	1015
962	Jiaxuan You, Tianyu Du, and Jure Leskovec. 2022.	<i>Proceedings of the 2019 27th ACM joint meeting on</i>	1016
963	Roland: graph learning framework for dynamic	<i>European software engineering conference and sym-</i>	1017
964	graphs. In <i>Proceedings of the 28th ACM SIGKDD</i>	<i>posium on the foundations of software engineering,</i>	1018
965	<i>conference on knowledge discovery and data mining,</i>	pages 807–817.	1019
966	pages 2358–2366.	Zeyang Zhang, Xin Wang, Ziwei Zhang, Zhou Qin,	1020
		Weigao Wen, Hui Xue, Haoyang Li, and Wenwu	1021
		Zhu. 2023b. Spectral invariant learning for dynamic	1022
		graphs under distribution shifts. <i>Advances in Neural</i>	1023
		<i>Information Processing Systems</i> , 36:6619–6633.	1024

1025 Qihang Zhao and Xiaodong Feng. 2022. Utilizing cita-
1026 tion network structure to predict paper citation counts:
1027 A deep learning approach. *Journal of Informetrics*,
1028 16(1):101235.

1029 Aoxiao Zhong, Dengyao Mo, Guiyang Liu, Jinbu Liu,
1030 Qingda Lu, Qi Zhou, Jiesheng Wu, Quanzheng Li,
1031 and Qingsong Wen. 2024. Logparser-llm: Advancing
1032 efficient log parsing with large language models. In
1033 *Proceedings of the 30th ACM SIGKDD Conference*
1034 *on Knowledge Discovery and Data Mining*, pages
1035 4559–4570.

1036 Jieming Zhu, Shilin He, Pinjia He, Jinyang Liu, and
1037 Michael R Lyu. 2023a. Loghub: A large collec-
1038 tion of system log datasets for ai-driven log analyt-
1039 ics. In *2023 IEEE 34th International Symposium*
1040 *on Software Reliability Engineering (ISSRE)*, pages
1041 355–366. IEEE.

1042 Yifan Zhu, Fangpeng Cong, Dan Zhang, Wenwen Gong,
1043 Qika Lin, Wenzheng Feng, Yuxiao Dong, and Jie
1044 Tang. 2023b. Wingnn: Dynamic graph neural net-
1045 works with random gradient aggregation window. In
1046 *Proceedings of the 29th ACM SIGKDD conference on*
1047 *knowledge discovery and data mining*, pages 3650–
1048 3662.

A Appendix

A.1 Details of datasets

Table 2 summarizes the basic setups of our datasets. The Arxiv dataset is publicly open to research purpose, where we anonymize the author name. HDFS is part of LogHub (Zhu et al., 2023a) which is freely open to research. Other two datasets are accessed from the Web without any privacy information.

A.2 Algorithms

Algorithm 1 and 2 illustrate the pipe of training and evaluation, respectively.

Algorithm 1: GraphLogDebugger: Online Training for Dynamic-Graph Anomaly Detection

Input : Training log
 $X_{\text{train}} = \{x_{t_0}, \dots, x_{t_{K-1}}\}$; GNN parameters θ ; text embedding model \mathcal{F} ; negative sampling ratio ρ ; threshold τ

Output : Trained parameters θ^* ; dynamic graph snapshot G_{t_K}

- 1 Initialize graph $\mathcal{G} = (\mathcal{V} = \emptyset, \mathcal{E} = \emptyset)$
- 2 **for** $k = 0, \dots, K - 1$ **do**
- 3 **Integrate incoming log entry**
- 4 Build subgraph g_k from x_{t_k} with object nodes \mathcal{V}_k^o , event nodes \mathcal{V}_k^e , and object–event links
- 5 $\mathcal{V} \leftarrow \mathcal{V} \cup \mathcal{V}_k^o$
- 6 $\mathcal{E}_k^+ \leftarrow$ object–event links in g_k
- 7 Sample \mathcal{E}_k^- with $|\mathcal{E}_k^-| = \rho \cdot |\mathcal{E}_k^+|$
- 8 **Embed nodes**
- 9 Compute object embeddings
 $h_o = f_\theta(\mathcal{G})$
- 10 Compute event embeddings
 $h_e = \mathcal{F}(\mathcal{V}_k^e)$
- 11 **Predict links and update parameters**
- 12 **foreach** $(o, e) \in \mathcal{E}_k^+ \cup \mathcal{E}_k^-$ **do**
- 13 | $s_{o,e} = \sigma(\text{MLP}(\text{reduce}(h_o, h_e)))$
- 14 Compute balanced BCE loss \mathcal{L}_k
- 15 Update parameters $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_k$
- 16 **Update dynamic graph**
- 17 $\mathcal{V} \leftarrow \mathcal{V} \cup \mathcal{V}_k^e$
- 18 Repeat for multiple epochs
- 19 **return** $\theta^* = \theta, G_{t_K} = \mathcal{G}$

Algorithm 2: GraphLogDebugger: Online Evaluation for Dynamic-Graph Anomaly Detection

Input : Test log
 $X_{\text{test}} = \{x_{t_K}, \dots, x_{t_{N-1}}\}$; trained GNN θ^* ; initial graph G_{t_K} ; text embedding model \mathcal{F} ; threshold τ

Output : Per-time link predictions
 $\{\mathcal{R}_{t_n}\}_{n=K}^{N-1}$; updated snapshot G_{t_N}

- 1 Initialize $\mathcal{G} \leftarrow G_{t_K}$
- 2 **for** $n = K, \dots, N - 1$ **do**
- 3 **Integrate incoming log entry**
- 4 Build subgraph g_n from x_{t_n} with objects \mathcal{V}_n^o , events \mathcal{V}_n^e , and observed links \mathcal{E}_n^+
- 5 $\mathcal{V} \leftarrow \mathcal{V} \cup \mathcal{V}_n^o$
- 6 **Embed nodes**
- 7 Compute object embeddings
 $h_o = f_{\theta^*}(\mathcal{G})$
- 8 Compute event embeddings
 $h_e = \mathcal{F}(\mathcal{V}_n^e)$
- 9 **Predict links**
- 10 **foreach** $(o, e) \in \mathcal{E}_n^+$ **do**
- 11 | $s_{o,e} = \sigma(\text{MLP}(\text{reduce}(h_o, h_e)))$
- 12 | $\hat{\ell}_{o,e} = \mathbb{1}[s_{o,e} \geq \tau]$
- 13 $\mathcal{R}_{t_n} \leftarrow \{(o, e, s_{o,e}, \hat{\ell}_{o,e})\}$
- 14 **Update dynamic graph**
- 15 $\hat{\mathcal{E}}_n^+ = \{(o, e) \in \mathcal{E}_n^+ \mid \hat{\ell}_{o,e} = 1\}$
- 16 $\hat{\mathcal{V}}_n^e = \{e \in \mathcal{V}_n^e \mid \exists o : (o, e) \in \hat{\mathcal{E}}_n^+\}$
- 17 $\mathcal{G} \leftarrow (\mathcal{V} \cup \hat{\mathcal{V}}_n^e, \mathcal{E} \cup \hat{\mathcal{E}}_n^+)$
- 18 **return** $\{\mathcal{R}_{t_n}\}_{n=K}^{N-1}, G_{t_N} = \mathcal{G}$

A.3 Model Design Space

We compare two variants in our experiments: (i) Plain (ungated) GAT. We first concatenate the entity-type and entity-ID embeddings and pass them through a feed-forward projection to obtain the initial representation e_0 . We then run multi-layer, multi-head GATConv on an entity–entity graph induced by shared content to propagate messages and obtain e_{GAT} , which we use as the final entity representation. (ii) Gated fusion. Starting from the same e_0 and e_{GAT} , we introduce a global learnable scalar gate α and adaptively combine them via a sigmoid: $e = (1 - \sigma(\alpha))e_0 + \sigma(\alpha)e_{\text{GAT}}$. This biases toward e_0 when the given signal is weak (or absent) and toward e_{GAT} when the signal is strong. Both variants share the same link-prediction head: we take the entity representation and the content representation (text and time embeddings concate-

Table 2: Statistics and details of the four datasets for tabular log debugging.

Dataset	Domain	#Entries	#Objects	Event Attr.	Obj Attr.	Anomaly Type
Arxiv	Sci. Pub.	20,000	17316	title	authors	Event/Object
HDFS	System Log	20,000	2150	Content	Component,EventId,BlockId	Object
Analyst	Finance	20,000	3901	headline	publisher	Event
Landslide	Geology	20,000	8565	description	title, category,trigger,country	Object

nated and then projected), compute their element-wise difference, and feed it to an MLP to output the link probability. We use A6000 GPUs to train the model.

A.4 Additional Visualization

Figure 3 visualizes the distribution of anomaly likelihood scores of our five evaluation tasks. The score distribution corroborates the main result in Table 1, that Analyst, Arxiv (Node), and HDFS are three tasks relatively easy, with the score distribution of anomalies and normal examples separate clearly. By contrast, the score of anomalies and normal examples mix up in Arxiv (Edge) and Landslide, indicating that these datasets are more difficult.

A.5 Notation

Table 3 lists all notations used in the paper.

A.6 AI Assistants In Research Or Writing

We use ChatGPT to polish our introduction (Section 1) and generate the notation table (Table 3), both of which have been checked manually. We also use ChatGPT to retrieve related works in the tabular log processing part by searching machine-learning based log processing methods.

A.7 Prompts in RAG

We list the prompt we used in our RAG baseline as follows:

```

1 context = """You are an expert at
2   analyzing author-paper relationships
3   in academic research.
4 DATASET CONTEXT: This is a dataset of
5   academic papers with their authors
6   and titles.
7   - Entities (authors): Research authors
8     who wrote the papers
9   - Content (titles): The titles of the
10  academic papers
11 - Edge: A connection between an author
12   and a paper title (indicating the
13   author contributed to that paper)
14 TASK: Determine if the specific author-
15   paper connection (edge) should exist
16   based on historical patterns.
17 EDGE ANALYSIS TARGET:

```

```

11 """
12
13 context += f"Author: {entity_name}\n"
14 context += f"Paper Title: {content_name
15   }\n\n"
16
17 if similar_contents:
18   context += "SIMILAR PAPERS AND THEIR
19     AUTHORS (for reference):\n"
20   context += "Use these examples to
21     understand what types of authors
22     typically work on similar
23     papers.\n\n"
24
25   for i, content_record in enumerate(
26     similar_contents[:10]):
27     content = content_record.get('
28       content', '')
29     entities = content_record.get('
30       related_entities', [])
31     similarity = content_record.get(
32       'similarity', 0.0)
33     num_records = content_record.get(
34       'num_records', 0)
35
36     context += f"{i+1}. Paper Title:
37       {content} (Similarity: {
38         similarity:.3f}, {
39         num_records} records)\n"
40     context += f"  Authors who
41       worked on this paper: {'', '
42       .join(entities) if entities
43       else 'None'}\n\n"
44 else:
45   context += "No similar papers found
46     in historical data.\n\n"
47
48 context += """ANALYSIS QUESTION:
49 Based on the similar papers and their
50 author patterns, should the
51 specified author-paper connection
52 exist?
53
54 EVALUATION CRITERIA:
55 1. Research Domain Match: Does the
56   author's expertise align with the
57   paper's topic?
58 2. Historical Patterns: Do authors with
59   similar expertise appear in similar
60   papers?
61 3. Authorship Likelihood: Is it
62   reasonable that this author would
63   contribute to this type of research?
64 4. Anomaly Detection: Does this
65   connection seem unusual or out of
66   place compared to patterns in
67   similar papers?
68
69 DECISION GUIDELINES:
70 - edge_exists = True: The author-paper

```

1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184

Table 3: **Notation**

Symbol	Type	Meaning
Tabular-log basics (Sec. 3)		
$X = \{x_0, \dots, x_{N-1}\}$	sequence	Time-ordered tabular log (entries).
x_n	entry	The n -th log entry.
$t_0 < \dots < t_{N-1}$	timestamps	Arrival times of entries.
x_n^m	attribute value	The m -th attribute in entry x_n .
M	integer	Number of attributes per entry.
$\{o_n^0, \dots, o_n^{P-1}\}$	set	Object attributes extracted from x_n .
P	integer	Number of object attributes in x_n ($P < M$).
s_n	text / node	Event attribute (one per entry; possibly text).
Q	integer	Number of feature attributes in x_n ($Q < M$).
y_n	label	Event anomaly label for x_n (0 normal, 1 abnormal).
y_n^p	label	Object anomaly label for object o_n^p (0/1).
Graphs and dynamics (Sec. 4.1–4.2)		
\mathcal{G}	dynamic graph	Evolving heterogeneous graph over time.
G_n	snapshot	Graph snapshot at time t_n (before merging g_n).
g_n	subgraph	Subgraph constructed from new entry x_n .
$G_{n+1} \setminus G_n$	graph diff	Increment between snapshots; here equal to g_n .
\mathcal{V}, \mathcal{E}	sets	Node and edge sets of the current graph.
v, e	node, edge	A node or an edge (generic).
\mathcal{V}_n^o	node set	Object nodes appearing in x_n .
\mathcal{V}_n^e	node set	New event nodes introduced by x_n (events are unique).
\mathcal{E}_k^+	edge set	Positive (observed) object–event links in g_k .
\mathcal{E}_k^-	edge set	Negative samples (non-existent object–event pairs).
$\hat{\mathcal{E}}_n^+$	edge set	Accepted/predicted-positive links at t_n .
$\hat{\mathcal{V}}_n^e$	node set	Accepted new events incident to $\hat{\mathcal{E}}_n^+$.
\mathcal{R}_{t_n}	set	Per-time link predictions/results at t_n .
$\{G_n\}_{n=0}^{N-1}$	sequence	The sequence of snapshots defining \mathcal{G} .
Modeling (GNN and scoring; Sec. 4.3)		
θ	parameters	Trainable parameters of the GNN.
$f_\theta(\cdot)$	mapping	GNN that computes object-node embeddings on \mathcal{G} .
\mathcal{F}	encoder	Text (and time-aware) embedding model for events.
h_o, h_e	vectors	Object and event embeddings, respectively.
$\text{reduce}(\cdot, \cdot)$	operator	Embedding combiner (e.g., concat/diff/dot).
$\text{MLP}(\cdot)$	mapping	Multi-layer perceptron used for scoring.
$\sigma(\cdot)$	function	Sigmoid activation.
$s_{o,e}$	score	Link-normality score for pair (o, e) .
$\hat{\ell}_{o,e}$	label	Predicted link label: $\mathbb{1}[s_{o,e} \geq \tau]$.
\mathcal{L}_k	loss	Balanced BCE loss at training step k .
η	scalar	Learning rate.
τ	threshold	Operating threshold for prediction.
ρ	ratio	Negative sampling ratio.
Data splits and indices		
$X_{\text{train}}, X_{\text{test}}$	sequences	Training and test splits (chronological).
K	integer	Index/time that separates train and test.
N	integer	Total number of entries/snapshots.
k, n	indices	Training step k , evaluation time n .
t_k, t_n	timestamps	Times associated with steps/entries.

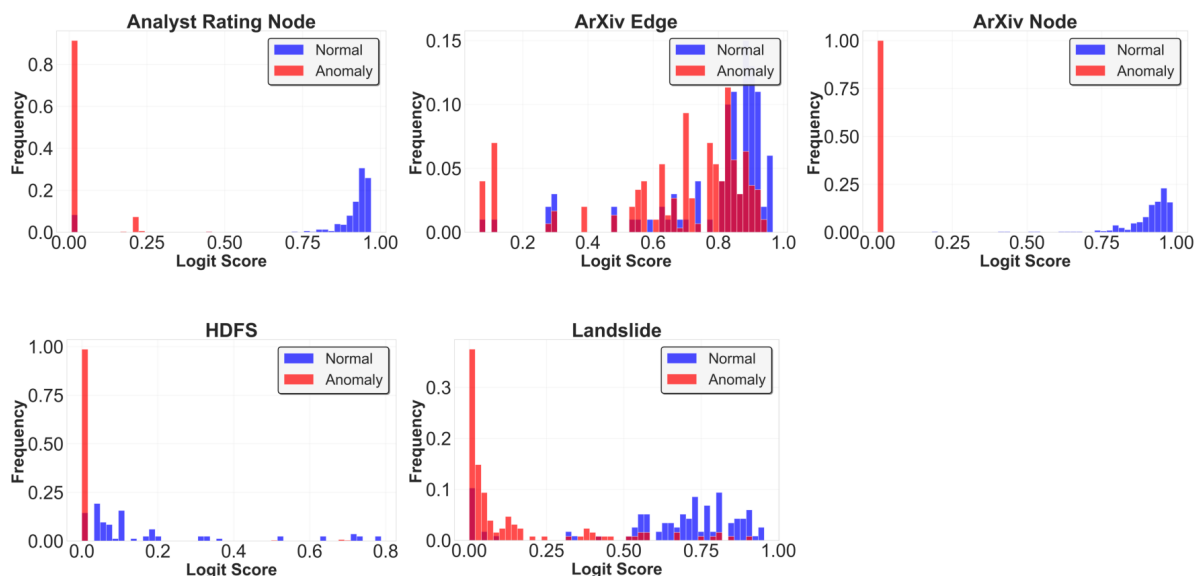


Figure 3: **Anomaly score distribution of five tasks by GraphLogDebugger.** Score distributions of anomalies and normal examples separate for simpler tasks and mix up for more difficult tasks.

```

1185 connection makes sense based on
1186 research area and historical
1187 patterns
1188 - edge_exists = False: The author seems
1189 misplaced or unlikely to work on
1190 this type of paper (anomalous edge)
1191 - Consider the research fields,
1192 methodologies, and typical author
1193 patterns shown in similar papers
1194 - An edge is anomalous if the author
1195 appears completely unrelated to the
1196 research domain of the paper
1197
1198 CONFIDENCE SCORING:
1199 - High confidence (0.8-1.0): Clear
1200 patterns in similar papers strongly
1201 support/reject the connection
1202 - Medium confidence (0.5-0.7): Some
1203 evidence but less certain
1204 - Low confidence (0.0-0.4): Limited
1205 historical data or unclear patterns
1206 """

```

Listing 1: Prompt: Arxiv

```

1208
1209
1210 context = """You are an expert at
1211 analyzing BlockId-log relationships
1212 in HDFS distributed file system logs
1213 .
1214
1215 DATASET CONTEXT: This is a dataset of
1216 HDFS system logs with their Block
1217 IDs and log contents.
1218 - Primary Focus: Block IDs (e.g.,
1219 blk_8215417782549978040,
1220 blk_161475555609545016) - unique
1221 identifiers for HDFS data blocks
1222 - Content (logs): The actual log
1223 messages and operations in the HDFS
1224 system that involve specific blocks

```

```

6 - Edge: A connection between a Block ID
7 and a log message (indicating the
8 block is involved in that log
9 operation)
10
11 SPECIAL NOTE: For HDFS anomaly detection
12 , we focus specifically on Block ID
13 connections to log messages.
14 Block IDs should appear BOTH in the
15 BlockId column AND within the log
16 content itself.
17
18 TASK: Determine if the specific Block ID
19 -log connection (edge) should exist
20 based on historical patterns.
21
22 EDGE ANALYSIS TARGET:
23 """
24
25 context += f"Block ID: {entity_name}\n"
26 context += f"Log Content: {content_name
27 }\n"
28 context += f"Content Analysis: Does '{
29 entity_name}' appear in the log
30 content? {'YES' if entity_name in
31 content_name else 'NO'}\n\n"
32
33 if similar_contents:
34 context += "SIMILAR LOG MESSAGES AND
35 THEIR BLOCK IDs (for reference)
36 :\n"
37 context += "Use these examples to
38 understand what types of Block
39 IDs typically appear in similar
40 log messages.\n\n"
41
42 for i, content_record in enumerate(
43 similar_contents[:10]):
44 content = content_record.get('
45 content', '')
46 entities = content_record.get('

```

1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264

```

1265         related_entities', [])
1266 27 similarity = content_record.get(
1267     'similarity', 0.0)
1268 28 num_records = content_record.get
1269     ('num_records', 0)
1270 29
1271 30 block_ids = [e for e in entities
1272     if e.startswith('blk_')]
1273 31 other_entities = [e for e in
1274     entities if not e.startswith
1275     ('blk_')]
1276 32
1277 33 context += f"{i+1}. Log Content:
1278     {content} (Similarity: {
1279     similarity:.3f}), {
1280     num_records} records)\n"
1281 34 context += f"    Block IDs in
1282     this log: {'', '.join(
1283     block_ids) if block_ids else
1284     'None'}\n"
1285 35 if other_entities:
1286 36     context += f"    Other
1287         entities: {'', '.join(
1288         other_entities[:3])
1289         }{'...' if len(
1290         other_entities) > 3 else
1291         ''}\n"
1292 37     context += "\n"
1293 38 else:
1294 39     context += "No similar log messages
1295         found in historical data.\n\n"
1296 40
1297 41 context += """ANALYSIS QUESTION:
1298     Based on the similar log messages and
1299     their Block ID patterns, should the
1300     specified Block ID-log connection
1301     exist?
1302 43
1303 44 EVALUATION CRITERIA:
1304 45 1. Block ID Presence: Does the Block ID
1305     appear within the log content itself
1306     ? (This is crucial for HDFS)
1307 46 2. Log Operation Match: Does the Block
1308     ID relate to the HDFS operation
1309     described in the log?
1310 47 3. Historical Patterns: Do similar Block
1311     IDs appear in similar log messages?
1312 48 4. HDFS Block Behavior: Is it reasonable
1313     that this Block ID would be
1314     involved in this type of operation?
1315 49 5. Content Consistency: Block ID should
1316     be consistent between the BlockID
1317     column and the log content
1318 50
1319 51 DECISION GUIDELINES:
1320 52 - edge_exists = True: The Block ID-log
1321     connection makes sense based on HDFS
1322     block operations and historical
1323     patterns
1324 53 - edge_exists = False: The Block ID
1325     seems unrelated to this log message
1326     (anomalous edge)
1327 54 - CRITICAL: If the Block ID does NOT
1328     appear in the log content, this is
1329     likely anomalous
1330 55 - Consider HDFS block operations like
1331     allocation, storage, replication
1332     shown in similar messages
1333 56 - An edge is anomalous if the Block ID
1334     appears completely unrelated to the

```

```

57     log operation
58 CONFIDENCE SCORING:
59 - High confidence (0.8-1.0): Clear Block
60   ID patterns and content consistency
61   strongly support/reject the
62   connection
63 - Medium confidence (0.5-0.7): Some
64   evidence but less certain about
65   Block ID relevance
66 - Low confidence (0.0-0.4): Limited
67   historical data or unclear Block ID
68   patterns
69 IMPORTANT: Focus specifically on Block
70 ID relationships - Components and
71 Event IDs are secondary for this
72 analysis.
73 """

```

```

1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353

```

Listing 2: Prompt: HDFS

```

1 context = f"""You are an expert at
2   analyzing entity-content
3   relationships.
4 EDGE ANALYSIS TARGET:
5 Entity: {entity_name}
6 Content: {content_name}
7 TASK: Determine if this entity-content
8   connection should exist based on
9   historical patterns.
10 """
11 if similar_contents:
12     context += "\nSIMILAR EXAMPLES:\n"
13     for i, content_record in enumerate(
14         similar_contents[:5]):
15         content = content_record.get('
16             content', '')
17         entities = content_record.get('
18             related_entities', [])
19         context += f"{i+1}. Content: {
20             content}\n    Related
21             entities: {'', '.join(
22             entities)}\n\n"
23
24 context += """
25 DECISION: Should this entity-content
26 connection exist?
27 - edge_exists = True: The connection
28   makes sense based on patterns
29 - edge_exists = False: The connection
30   seems anomalous
31 """

```

```

1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389

```

Listing 3: Prompt: Analyst and Landslide