

---

000 TOWARDS GENERALIZABLE CONTEXT-AWARE  
001 ANOMALY DETECTION: A LARGE-SCALE BENCH-  
002 MARK IN CLOUD ENVIRONMENTS  
003  
004  
005

006 **Anonymous authors**

007 Paper under double-blind review  
008  
009

010  
011 ABSTRACT  
012

013 Anomaly detection in cloud environments remains both critical and challenging.  
014 Existing context-level benchmarks typically focus on either metrics or logs and  
015 often lack reliable annotation, while most detection methods emphasize point  
016 anomalies within a single modality, overlooking contextual signals and limiting  
017 real-world applicability. Constructing a benchmark for context anomalies that  
018 combines metrics and logs is inherently difficult: reproducing anomalous sce-  
019 narios on real servers is often infeasible or potentially harmful, while generating  
020 synthetic data introduces the additional challenge of maintaining cross-modal con-  
021 sistency. We introduce CLOUDANOBENCH, a large-scale benchmark for context  
022 anomalies in cloud environments, comprising 28 anomalous scenarios and 16 de-  
023 ceptive normal scenarios, with 1,252 labeled cases and roughly 200,000 log and  
024 metric entries. Compared with prior benchmarks, CLOUDANOBENCH exhibits  
025 higher ambiguity and greater difficulty, on which both prior machine learning  
026 methods and vanilla LLM prompting perform poorly. To demonstrate its utili-  
027 ty, we further propose CLOUDANOAGENT, an LLM-based agent enhanced by  
028 symbolic verification that integrates metrics and logs. This agent system achieves  
029 substantial improvements in both anomaly detection and scenario identification  
030 on CLOUDANOBENCH, and shows strong generalization to existing datasets. To-  
031 gether, CLOUDANOBENCH and CLOUDANOAGENT lay the groundwork for ad-  
032 vancing context-aware anomaly detection in cloud systems.  
033

034 1 INTRODUCTION  
035

036 Ensuring the stability and availability of large-scale cloud systems is of great importance  
037 (Kazemzadeh & Jacobsen, 2009; Bu et al., 2018; Zhang et al., 2015). Accurate detection meth-  
038 ods that can also identify among anomaly scenarios are essential to mitigate potential losses (Zhang  
039 et al., 2018; Barbhuiya et al., 2018a). Large-scale cloud systems usually generate abundant logs and  
040 expose various metrics, both of which serve as some of the most valuable data sources for anomaly  
041 detection (Lin et al., 2016; Nandi et al., 2016).  
042

043 Numerous benchmarks have been proposed for cloud anomaly detection such as (Oliner & Stearley,  
044 2007; Xu et al., 2009; Akmeemana et al., 2025). However, most existing research and benchmarks  
045 for cloud anomaly detection have focused on point anomalies, where deviations are identified in iso-  
046 lation within a single modality, such as metrics or logs. Although these benchmarks have provided  
047 the community with relevant evaluation testbeds, they capture only a narrow slice of the anomaly  
048 landscape and often fail to reflect the complexity of real cloud environments. In practice, anoma-  
049 lies emerge from the interplay between multiple system components and are best understood in the  
050 context of both logs and metrics. Context anomalies fill this gap by modeling dependencies across  
051 modalities and incorporating richer scenario-level information. These anomalies account for more  
052 realistic cases such as failures caused by combined metric trends and log events. Existing context  
053 anomaly data sets (Makanju et al., 2009; Islam et al., 2025) focus on a singular modality that cap-  
tures either metrics or logs and lacks clear annotation. This underlines the importance of novel  
improved contextual anomaly benchmarks that represent real-world operational demands.

054 Despite the practical relevance of contextual anomalies in cloud environments, constructing a bench-  
055 mark that faithfully represents these anomalies is a challenging task. Many critical safety anomaly  
056 scenarios **cannot be replicated safely without irreversibly damaging the infrastructure**. Ad-  
057 ditionally, synthetic anomaly data needs to be generated carefully to **capture the nuanced rela-**  
058 **tionships between metrics and logs** while being representative of real world incidences. These  
059 difficulties highlight why existing datasets have fallen short and why creating a large-scale, multi-  
060 modal benchmark remains a fundamental challenge.

061 To address these challenges, we introduce CLOUDANOBENCH<sup>1</sup>, a large-scale benchmark for  
062 context-aware anomaly detection in cloud environments. Unlike existing benchmarks, it jointly  
063 incorporates both metrics and logs, comprising 1,252 scenario labeled cases across 28 anomalous  
064 scenarios and 16 deceptive normal scenarios (approximately 200K lines), with explicit anomaly and  
065 scenario annotations. Moreover, we design deceptive normal scenarios, where anomalous-looking  
066 metric patterns are explained by benign log events. These properties make CLOUDANOBENCH  
067 substantially more ambiguous and challenging than prior datasets.

068 To illustrate the utility of CLOUDANOBENCH, we propose CLOUDANOAGENT, an LLM-based  
069 agent framework enhanced with symbolic verification that jointly leverages both metrics and logs.  
070 Traditional approaches are often limited in handling unstructured log data, which contains criti-  
071 cal temporal and behavioral context (Abdallah et al., 2024; Lou et al., 2019). CLOUDANOAGENT  
072 integrates a Fast and Slow Detection mechanism, leveraging the reasoning capabilities of LLMs  
073 to analyze metrics alongside contextual event logs, thereby ensuring both responsiveness and ac-  
074 curacy in detection. Furthermore, we design a symbolic verification module (Sun & Bookman,  
075 1994), which performs statistical validation over metric patterns and regex-based matching over log  
076 events for each anomaly scenario. Acting as a critic to the *Integrated Agent*, this component fur-  
077 ther improves detection performance and reduces false positives. As a result, CLOUDANOAGENT  
078 demonstrates strong performance in both anomaly detection and scenario identification, while also  
079 exhibiting generalizability beyond the capabilities of existing methods.

080 We evaluate CLOUDANOBENCH and CLOUDANOAGENT through extensive experiments.  
081 CLOUDANOAGENT achieves up to nearly 20% higher F1-score on anomaly detection and 15%  
082 on scenario identification over baselines (e.g., machine learning methods and vanilla LLM prompt-  
083 ing), which often struggle with deceptive normal cases and accurate scenario identification. Notably,  
084 CLOUDANOAGENT also shows competitive performance on log-only, point-anomaly datasets, high-  
085 lighting its strong generalizability beyond its original design for context-level multimodal anomalies.

086 Our contributions can be summarized as:

- 087 • We introduce CLOUDANOBENCH, a large-scale benchmark for context anomalies that jointly  
088 includes metrics and logs, designed to evaluate both anomaly detection and scenario identification.
- 089 • We propose CLOUDANOAGENT, an LLM-based agent framework guided by symbolic verifica-  
090 tion, which simultaneously leverages metrics and logs and exhibits strong generalization ability.
- 091 • Experimental results demonstrate the challenging and deceptive nature of CLOUDANOBENCH, as  
092 well as the superior performance of CLOUDANOAGENT in both anomaly detection and scenario  
093 identification, achieving higher performance than compared methods and strong generalization on  
094 other datasets.

## 096 2 RELATED WORK

098 Anomaly detection is the process of identifying and extracting unexpected behaviors and pat-  
099 terns from the data. As cloud systems continue to evolve and grow rapidly, detecting anomalies  
100 has become critical to ensure cloud stability and reliability. Early deep learning approaches like  
101 LogAnomaly (Meng et al., 2019) used LSTMs to model normal sequences of log keys and detect  
102 deviations. This was advanced by LogBERT (Guo et al., 2021), which adapted Transformer to  
103 learn deep contextual representations of logs through self-supervised pre-training. More recently,  
104 LogLLM (Guan et al., 2024) has demonstrated the power of LLMs, coordinating BERT and Llama  
105 to extract semantic vectors and classify log sequences. These methods are designed specifically  
106 for log-based point anomalies. Therefore, we use them as baselines for evaluation in this work to  
107 demonstrate CLOUDANOAGENT’s performance when the data is limited to log only.

<sup>1</sup>The current benchmark is available at <https://anonymous.4open.science/r/CloudAnoBench-17ED>

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

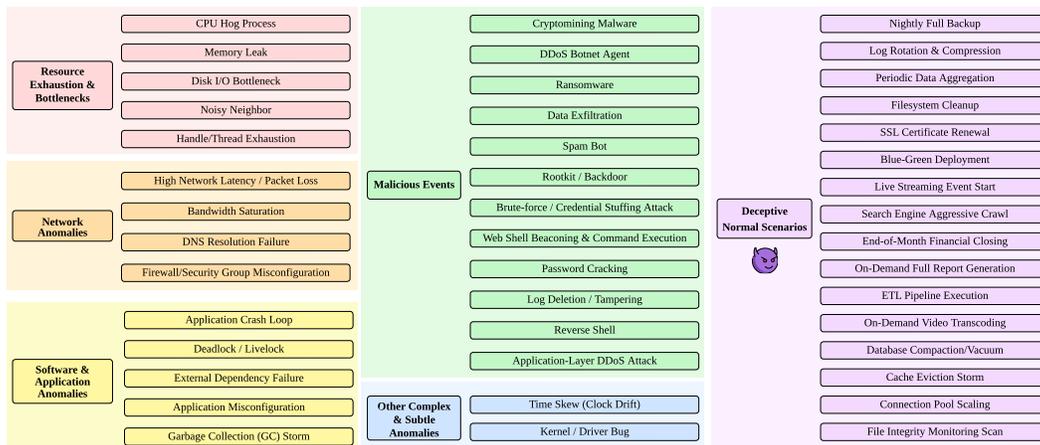


Figure 1: Overview Taxonomy of CloudAnoBench

TraceBench (Zhou et al., 2014) is a dataset of trace anomalies collected from real-world distributed systems, containing 370,000 traces and covering 17 types of anomalies. HPC dataset (Makanju et al., 2009) offers log-based context anomalies collected from a high performance computing cluster. However, it does not contain metrics data and its data are not labeled, limiting its applicability for evaluating techniques. IBM Cloud Console (Islam et al., 2025) provides context anomaly data from cloud systems that are properly labeled. However, this dataset contains metrics data only. More recently, LO2 (Bakhtin et al., 2025) is a novel comprehensive dataset of logs, metrics, and traces from a microservice system, encompassing multiple data modalities. Although not specific to cloud environments, it highlights the significance of using contextual data for effective anomaly detection.

### 3 CLOUDANOBENCH

| Dataset                                | Metrics & Logs | Type    | Labeled       | # Scenarios |
|--|----------------|---------|---------------|-------------|
| BGL (Oliner & Stearley, 2007)          | ✗              | Point   | ✓             | –           |
| Thunderbird (Oliner & Stearley, 2007)  | ✗              | Point   | ✓             | –           |
| HPC (Makanju et al., 2009)             | ✗              | Context | ✗             | –           |
| HDFS_v1 (Xu et al., 2009)              | ✗              | Point   | ✓             | 17          |
| BETH (Highnam et al., 2021)            | ✗              | Point   | ✓             | –           |
| IBM Cloud Console (Islam et al., 2025) | ✗              | Context | ✓             | –           |
| RS-Anomic (Akmeemana et al., 2025)     | ✗              | Point   | ✓             | 10          |
| <b>CloudAnoBench (Ours)</b>            | ✓              | Context | ✓ (scenarios) | 28          |

Table 1: Comparison of CloudAnoBench with existing datasets, highlighting its various scenarios, clear annotations with scenarios, and coverage of context anomalies that jointly include both metrics and logs.

We introduce CLOUDANOBENCH, a large-scale benchmark for evaluating context-aware anomaly detection in cloud environments. CLOUDANOBENCH consists of 1,252 cases covering 28 anomalous scenarios and 16 normal scenarios, as presented in Figure 1, totaling approximately 200,000 lines of data. Unlike existing benchmarks, shown in Table 1, it jointly includes context-level metrics and log data, with each case annotated by explicit anomaly and scenario labels. Moreover, we emphasize that anomaly detection methods should not raise alerts under normal operating conditions, as false alarms can incur significant financial and operational costs. To this end, CLOUDANOBENCH also incorporates a diverse set of normal scenarios that occur frequently in practice yet remain highly challenging for detectors. These properties make CLOUDANOBENCH substantially more challenging for prior machine learning and LLM prompting approaches, while also covering a broader spectrum of anomaly scenarios and aligning more closely with real-world cloud operations. More details

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

---

**Algorithm 1:** CLOUDANOBENCH Generation with Manual Review

---

**Input:** Scenario set  $\mathcal{S}$  with anomaly and normal cases, example  $e$   
**Output:** Final dataset  $\mathcal{D}$  with paired metrics.csv and log files

```

1 foreach scenario  $s \in \mathcal{S}$  do
2   Action: Construct structured prompt  $P_s$ 
3   if example  $e$  exists for  $s$  then
4     | Attach  $e$  as one-shot reference to guide consistency in format and content;
5   end
6   for  $i \leftarrow 1$  to number of required cases do
7     (Generate Metrics Data)
8      $f_{\text{metrics}} \leftarrow \text{LLM.GenerateCode}(P_s)$  with code_execution;
9     metrics.csv  $\leftarrow \text{Execute}(f_{\text{metrics}})$ ;
10    Action: anomaly pattern and data format verification
11    (Generate Log Data)
12    log  $\leftarrow \text{LLM.GenerateLog}(P_s, \text{aligned with metrics.csv})$ ;
13    Action: alignment with metrics, avoid metric description, insert benign noise
14    (Verification and Manual Review)
15    if  $\text{GPT-4o.Verifier}(\text{metrics}, \text{log}) = \text{pass}$  and  $\text{ManualReview} = \text{pass}$  then
16      | Append  $(\text{metrics}, \text{log}, \text{label})$  to  $\mathcal{D}$ ;
17    end
18    else
19      | Request regeneration of the case;
20    end
21  end
22 end
23 return  $\mathcal{D}$ ;

```

---

of CLOUDANOBENCH are provided in Appendix D and Appendix E, while additional information about the compared datasets is presented in Appendix G.

### 3.1 SCENARIO EXTRACTION

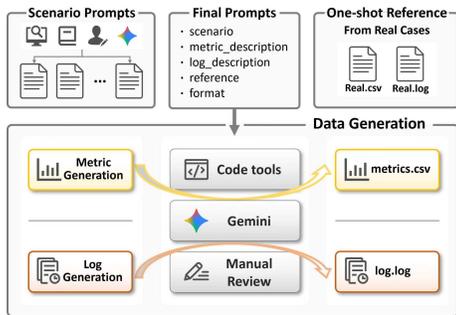


Figure 2: CloudAnoBench Construction

(if applicable), description, and corresponding metrics and logs. To ensure clarity and consistency across the dataset, instead of retaining the original heterogeneous formats, we employ Gemini-2.5-Flash to standardize every case into a structured prompt format, thereby facilitating reproducibility and enabling further data generation.

### 3.2 DATA GENERATION

To construct CLOUDANOBENCH, we leverage Gemini-2.5-Pro with scenario-specific prompts and tool invocation. Since many anomalies (e.g., cryptomining) cannot be safely reproduced on real

---

216 servers, and prior studies have demonstrated the feasibility of LLM-based dataset construction, we  
217 adopt an LLM-assisted pipeline. One real-world examples are provided as one-shot references to  
218 stabilize generation and ensure consistency. Each case consists of synchronized `metrics.csv`  
219 and `log` files that simulate realistic cloud environments. The full process, including prompt con-  
220 struction, metric/log generation, verification, and manual review, is summarized in Algorithm 1.  
221

222 **Metrics Data.** Each case contains multivariate system metrics (CPU, GPU, memory, disk I/O, and  
223 network), spanning 90 lines over 450 seconds. Five canonical anomaly patterns are simulated (spike,  
224 dip, gradual increase, gradual decrease, fluctuation), with values sampled under realistic constraints  
225 to ensure plausibility. Metric data is generated via LLM with `code_execution`, which produces  
226 valid `.csv` files.

227 **Log Data.** Log data spans the same 450-second window, with 40–80 entries temporally aligned to  
228 the metrics. Logs are generated to avoid direct metric restatements while incorporating benign noise  
229 (e.g., SSH logins, scheduled tasks), thereby reflecting real-world conditions.  
230

### 231 3.3 MANUAL REVIEW WITH LLM 232

233 To ensure data quality, we combine automatic checks with GPT-4o verification and human review.  
234 GPT-4o filters low-quality or inconsistent cases, while human reviewer adjudicate flagged samples  
235 and confirm anomaly classifications, ensuring realism and consistency across modalities.

236 As a result, CLOUDANOAGENT offers a large-scale, systematically validated, and multimodal  
237 benchmark that faithfully reflects realistic cloud anomalies while maintaining consistency, diver-  
238 sity, and reproducibility. Examples can be found in Appendix F. [To verify that CLOUDANOAGENT  
239 reflects real-world data patterns, we conducted an additional comparison experiment using the RS-  
240 Anomic dataset \(Akmeemana et al., 2025\). RS-Anomic provides multivariate metrics collected from  
241 the RobotShop Instana microservice system and is widely used for anomaly detection and root-cause  
242 analysis in production-like settings. The experiment evaluates the similarity between the temporal  
243 behaviors in CLOUDANOAGENT and those in RS-Anomic. Full results are reported in Appendix C.](#)  
244

## 245 4 CLOUDANOAGENT 246

247 In this work, we propose CLOUDANOAGENT, an LLM-based agent guided by symbolic verification  
248 for cloud anomaly detection. **The main innovation lies in proposing a novel context engineering  
249 framework specifically for cloud anomaly detection.** Specifically, Our framework is comprised  
250 of two main modules: Fast and Slow Detection and Symbolic Verifier. Fast and Slow Detection  
251 consists of three agents: a Metrics Agent, a Log Agent, and an Integrated Agent that jointly produce  
252 an initial decision. This initial decision is then validated by a Symbolic Verifier module, which  
253 applies rigorous rule-based symbolic techniques to ensure reliability. This validation also helps to  
254 refine the hypotheses for the Integrated Agent, effectively creating a feedback loop between the  
255 two modules. This hybrid approach significantly enhances detection precision and reduces false  
256 positives.

### 257 4.1 FAST AND SLOW DETECTION 258

259 To capture both short-term anomalies (e.g., sudden spikes or drops) and context-dependent anom-  
260 alies (e.g., deviations explainable only through associated log events), CLOUDANOAGENT adopts  
261 a Fast and Slow Detection mechanism. This dual-phase strategy enables (1) responsiveness, by  
262 performing real-time detection of abrupt anomalous signals, (2) interpretability, by uncovering un-  
263 derlying causes and temporally extended anomaly patterns through deeper log reasoning, and (3)  
264 precision, by reducing false positives through joint reasoning. Both Fast and Slow Detection are  
265 performed by LLM-based agents, yet they diverge in terms of their invocation timing, agent compo-  
266 sition, and, most importantly, the design of their respective context engineering workflows.

267 **Fast Detection** is performed by a dedicated Metrics Agent, designed to continuously monitor real-  
268 time system signals within a sliding time window. This agent summarizes the behavior of metrics  
269 and outputs a structured report indicating whether an anomaly was detected and which metrics were  
affected, thereby providing the responsiveness needed for timely intervention.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

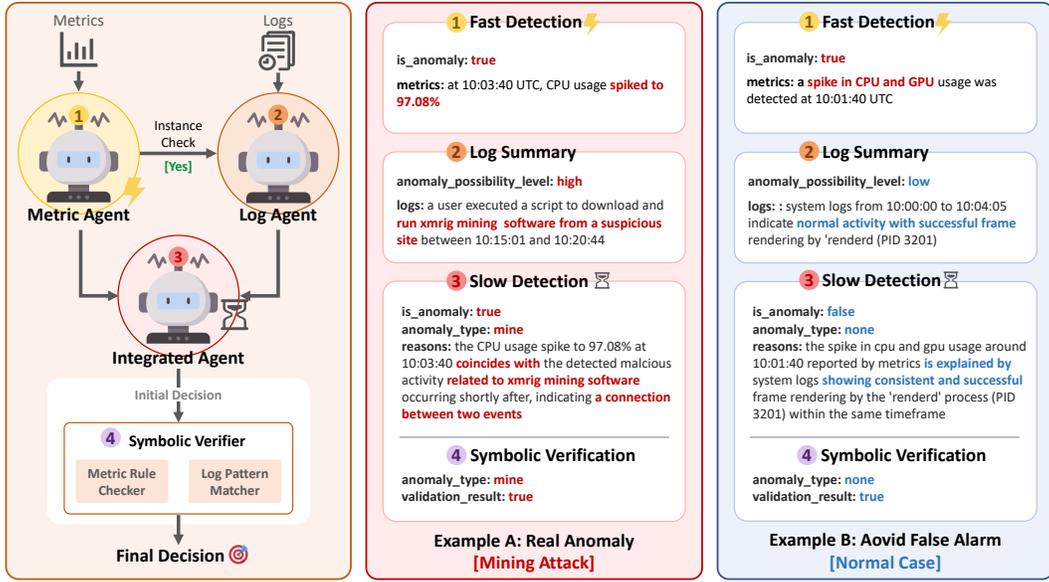


Figure 3: Overview of CloudAnoAgent. The red and blue boxes illustrate example outputs of CloudAnoAgent for different cases (simplified version).

**Slow Detection** is triggered in an event-driven manner upon anomaly signals raised by the Metrics Agent. Prior studies have shown that metric-based methods often suffer from high false positive rates in real-world cloud systems due to the lack of contextual information (Barbhuiya et al., 2018b; Kumari et al., 2024). To address this, our approach complements time-series signals with log-based reasoning to determine whether the detected anomaly reflects a truly abnormal event or a benign fluctuation. Beyond binary validation, this phase also identifies the likely root cause and higher-level anomaly event behind the observed deviations.

Specifically, a Log Agent processes log entries temporally aligned with the anomalous window, extracting the semantic meaning of system behaviors and assigning an anomaly likelihood level (i.e., low, medium, or high). These outputs, together with the response from the Metrics Agent, are fed into the Integrated Agent, which integrates both perspectives to perform joint anomaly verification and fine-grained classification. It determines whether the anomaly is valid and, if so, outputs the corresponding anomaly scenario along with a natural language explanation. This completes the end-to-end reasoning cycle in CloudAnoAgent’s Fast and Slow Detection mechanism.

## 4.2 SYMBOLIC VERIFIER

Despite the improved performance of our Fast and Slow Detection over other methods in terms of accuracy and false positive rate, the inherent stochasticity of LLM outputs raises concerns about reliability. The symbolic verifier is introduced because even with the fast-slow detection pipeline, LLM outputs can still exhibit stochastic errors, leading to occasional misclassification in highly deceptive cases. Its role is to provide a stable layer of constraint checking that prevents major category-level mistakes rather than determining fine-grained scenarios.

To address this, we introduce a symbolic verifier that not only validates the Integrated Agent’s initial decisions but also helps the system reflect and refine them when inconsistencies arise.

$$\text{Verifier}(M, L, s) = \begin{cases} 1, & \text{if } \text{VerifyMetric}(M, s) = 1 \wedge \text{VerifyLog}(L, s) = 1, \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

Specifically, for each predefined anomaly scenario  $s \in \mathcal{S}$ , the symbolic verifier applies two complementary checks: (1) a metric rule checker, which evaluates whether the observed metric sequence  $M$  satisfies the expected statistical signature of  $s$ , and (2) a log pattern matcher, which uses regex-based

---

keyword matching to confirm whether the log entries  $L$  contain events consistent with  $s$ . Formally, let  $\text{VerifyMetric}(M, s) \in \{0, 1\}$  denote the binary outcome of metric validation for scenario  $s$ , and  $\text{VerifyLog}(L, s) \in \{0, 1\}$  the corresponding outcome for log validation. The symbolic verifier accepts scenario  $s$  if and only if both conditions hold:

For the detector’s predicted scenario  $\hat{s}$ , if  $\text{Verifier}(M, L, \hat{s}) = 1$ , the decision is confirmed; otherwise, the system iteratively evaluates all  $s \in \mathcal{S}$  and either converges to a consistent alternative or abstains after a maximum retry limit. Appendix J provides an simplified example for the mining scenario. For rule construction, we use static, non-learned rules: an LLM first generates candidate metric and log rules, and human experts then review them for correctness and generality. Crucially, reviewers do not see CloudAnoBench’s data distribution, ensuring the rules remain dataset-agnostic and avoid overfitting.

## 5 EXPERIMENTS

This section shows our experimental setup, including datasets, compared methods, evaluation metrics, and implementation details. We evaluate detection methods on precision, recall, false positive rate (FPR) and F1-score on anomaly detection and accuracy on scenario identification (SI).

### 5.1 EXPERIMENTAL SETTINGS

#### 5.1.1 DATASETS

We conducted extensive experiments on four datasets to evaluate the effectiveness of the proposed methods for anomaly detection in cloud environments. These datasets include **BGL**, **Thunderbird**, **HDFS\_v1**, and our proposed benchmark, **CLOUDANOAGENT**. Table 1 provides an overview of these datasets. Although the IBM Cloud Console dataset also contains context-level anomalies, it is excluded from our evaluation because its anomalies are labeled using IBM’s proprietary internal monitoring tools, making it overly enterprise-specific and less suitable for general benchmarking (Islam et al., 2021). The use of BGL, Thunderbird, and HDFS\_v1 further demonstrates that while **CLOUDANOAGENT** is primarily designed to address context-level anomalies involving both metrics and logs, it also achieves strong performance on point-level anomaly datasets.

#### 5.1.2 COMPARED METHODS

We compare our proposed **CLOUDANOAGENT** framework against three categories of baselines: machine learning (ML) methods, vanilla LLM prompting, and log-only detection methods.

- **Machine Learning Algorithms.** We include several widely used machine learning baselines, including tree-based models such as *Isolation Forest* and *Decision Tree*, linear models such as *Logistic Regression*, unsupervised clustering methods such as *K-Means* and the *Rarity Model*, as well as the *OOV Detector*, which flags infrequent or previously unseen log tokens. These methods rely on fixed statistical patterns or token rarity assumptions, provide efficient and lightweight detection, and have been applied to anomaly detection tasks in microservice APIs (Bakhtin et al., 2025).
- **Vanilla LLM Prompting** To isolate the effectiveness of CloudAnoAgent’s agent architecture and symbolic verification from the raw capabilities of LLMs, we introduce a simple LLM baseline. In this setting, both the metrics data and the corresponding log text are directly fed into the LLM, which is prompted to output a binary anomaly decision along with a categorical scenario identification, without any agent decomposition or symbolic verification.
- **Log-Only Detection Methods.** We also compare against log-specific anomaly detection approaches, including LogAnomaly (Meng et al., 2019), LogLLM (Guan et al., 2024), and LogBERT (Guo et al., 2021). These methods are designed primarily for point anomaly detection in logs and serve as strong baselines to evaluate how well **CLOUDANOAGENT** generalizes to log-only detection settings.

## 5.2 EVALUATION METRICS

In this paper, we evaluate the performance of the proposed methods on both **anomaly detection** and **scenario identification**. For anomaly detection, the task is modeled as a binary classification problem. We report Precision, Recall, FPR, and F1-score, which together capture the reliability of positive alerts, the ability to identify true anomalies, the tendency to produce false alarms, and the overall balance between precision and recall.

For scenario identification (SI), we formulate the task as a multi-class classification problem and evaluate performance using accuracy (Acc). Let  $N$  denote the total number of cases,  $y_i$  the ground-truth label of the  $i$ -th case, and  $\hat{y}_i$  the predicted label. Accuracy is defined as:

$$\text{Acc} = \frac{|\{i \mid \hat{y}_i = y_i\}|}{N}, \quad (2)$$

where  $|\{i \mid \hat{y}_i = y_i\}|$  denotes the number of correctly classified cases. In other words, accuracy directly measures the proportion of correctly identified anomaly scenarios in the benchmark.

## 5.3 IMPLEMENTATION DETAILS

For both CLOUDANOAGENT and the LLM baseline, we accessed LLMs via API calls using the default inference parameters provided by each vendor. [Since ML methods cannot directly process textual log information, we add a log parser so that they can make predictions using both metrics and logs.](#) Moreover, since the BGL, Thunderbird, and HDFS\_v1 datasets only contain log data, the CLOUDANOAGENT was evaluated on these datasets without the metrics agent component. All experiments were conducted on a computing cluster equipped with two NVIDIA A100 GPUs. Some example prompts are provided in Appendix I.

## 6 RESULTS

| Setting  | Model                 | Anomaly Detection |                  |                  |                  | SI               |
|--|-----------------------|-------------------|------------------|------------------|------------------|------------------|
|  |                       | P                 | R                | FPR              | F1               | Acc              |
| ML Algorithms                                      | IsolationForest       | 54.0              | 98.1             | 100.0            | 69.7             | 53.5             |
|  | KMeans                | 55.3              | 93.0             | 90.0             | 69.4             | 55.2             |
|  | RarityModel           | 76.5              | 32.8             | 12.1             | 45.9             | 57.9             |
|  | OOV Detector          | 60.4              | 63.2             | 49.6             | 61.8             | 57.4             |
|  | <b>Average</b>        | <b>61.55</b>      | <b>71.78</b>     | <b>63.0</b>      | <b>61.95</b>     | <b>55.0</b>      |
| Vanilla LLM  | Gemini 2.5 Flash      | 66.9              | 91.4             | 54.3             | 77.2             | 66.9             |
|  | Gemini 2.5 Flash Lite | 62.4              | 89.0             | 64.5             | 73.3             | 61.2             |
|  | Gemini 2.0 Flash Lite | 60.4              | 86.1             | 67.7             | 71.0             | 60.0             |
|  | GPT-4o Mini           | 58.9              | 83.6             | 69.9             | 69.1             | 68.8             |
|  | GPT-4o                | 66.7              | 88.4             | 52.9             | 76.1             | 58.1             |
|  | Qwen Plus             | 60.3              | 83.9             | 66.3             | 70.2             | 60.3             |
|  | <b>Average</b>        | <b>62.6</b>       | <b>87.1</b>      | <b>62.6</b>      | <b>72.8</b>      | <b>62.6</b>      |
| CloudAnoAgent<br>w/o and with<br>symbolic verifier | Gemini 2.5 Flash      | 86.9/90.0         | 94.3/96.6        | 17.0/12.8        | 90.4/93.2        | 73.8/78.6        |
|  | Gemini 2.5 Flash Lite | 85.9/88.7         | 91.9/92.2        | 18.1/14.1        | 88.8/90.5        | 73.2/76.3        |
|  | Gemini 2.0 Flash Lite | 81.9/87.9         | 90.8/90.5        | 24.1/14.9        | 86.1/89.2        | 70.7/74.5        |
|  | GPT-4o Mini           | 87.5/90.0         | 89.5/90.8        | 15.3/12.1        | 88.5/90.4        | 71.9/73.1        |
|  | GPT-4o                | 91.5/92.5         | 93.6/95.6        | 10.4/9.3         | 92.5/94.0        | 74.8/79.2        |
|  | Qwen Plus             | 87.1/88.7         | 92.2/92.7        | 16.3/14.2        | 89.6/90.6        | 74.5/76.3        |
|  | <b>Average</b>        | <b>86.8/89.6</b>  | <b>92.0/93.1</b> | <b>16.9/12.9</b> | <b>89.3/91.3</b> | <b>73.2/76.3</b> |

Table 2: Evaluation results on CloudAnoBench: Precision (P), Recall (R), False Positive Rate (FPR), and F1-score for anomaly detection, and Accuracy (Acc) for scenario identification (SI).

This section presents our results on both anomaly detection and anomaly scenario identification. We analyze how well different methods detect anomalies, challenges of identifying specific anomaly scenarios, and evaluate the generalization performance of CLOUDANOAGENT across diverse datasets and anomaly types.

## 6.1 ANOMALY DETECTION PERFORMANCE

Our evaluation on anomaly detection (Table 2) demonstrates that CLOUDANOAGENT consistently outperforms all baselines on CLOUDANOBENCH, achieving the highest average F1-score (91.3%) and the lowest average FPR (12.9%). These results highlight its ability to capture true anomalies while suppressing false alarms, even in deliberately misleading normal cases. In contrast, vanilla LLM prompting attains competitive Recall (87.1% on average) but suffers from extremely high FPR (62.6% on average), indicating poor robustness against deceptive normal patterns and raising concerns about practical deployment costs. ML algorithms such as KMeans and RarityModel achieve recall above 95%, but exhibit poor precision and FPR. Finally, incorporating symbolic verification further improves CLOUDANOAGENT, reducing FPR by 4% and increasing F1 by 2% on average, underscoring the verifier’s role as a reliable critic that enhances detection accuracy. [More ablation study results can be found in Appendix B.](#)

## 6.2 CHALLENGES IN SCENARIO IDENTIFICATION

| CloudAnoAgent Model   | Number of Scenarios |       |       |
|-----------------------|---------------------|-------|-------|
|                       | 5                   | 10    | 29    |
| Gemini 2.5 Flash      | 100                 | 98.25 | 78.62 |
| Gemini 2.5 Flash Lite | 100                 | 96.49 | 76.28 |
| Gemini 2.0 Flash Lite | 97.35               | 91.23 | 74.52 |

Table 3: Performance of CLOUDANOAGENT across different numbers of anomaly scenarios.

However, as shown in Table 3, as the number of anomaly scenarios increases, scenario identification accuracy decreases significantly. This is because more scenarios require longer context windows and introduce higher semantic similarity across scenarios, making fine-grained identification more challenging for LLM-based agents.

## 6.3 GENERALIZATION OF CLOUDANOAGENT

| Model                | HDFS_v1     |             |             | ThunderBird |             |             | BGL         |             |             |
|----------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                      | Prec.       | Rec.        | F1          | Prec.       | Rec.        | F1          | Prec.       | Rec.        | F1          |
| LogLLM               | 99.4        | 100.0       | 99.7        | 96.6        | 96.6        | 96.6        | 86.1        | 97.9        | 91.6        |
| LogAnomaly           | 96.0        | 94.0        | 95.0        | -           | -           | -           | 97.0        | 94.0        | 96.0        |
| LogBERT              | 87.0        | 78.1        | 82.3        | 96.8        | 96.5        | 96.6        | 89.4        | 92.3        | 90.8        |
| <b>CloudAnoAgent</b> | <b>89.5</b> | <b>95.1</b> | <b>92.2</b> | <b>89.6</b> | <b>96.4</b> | <b>92.9</b> | <b>98.5</b> | <b>99.7</b> | <b>99.1</b> |

Table 4: Performance comparison on three log-only datasets with point anomaly type.

It is worth noting that while CLOUDANOAGENT is primarily designed for context-level anomaly datasets that include both metrics and logs, it can still be adapted to different data modalities and anomaly types through simple modifications to the prompting rules and lightweight adjustments to the agent workflow. As shown in Table 4, CLOUDANOAGENT achieves competitive performance on HDFS\_v1, ThunderBird, and BGL, when compared to LogLLM, LogAnomaly, and LogBERT—methods that are specifically designed for log-based point anomaly detection. The results of these baseline methods on the three datasets are reported from each original paper, and implementation details can be found in Appendix H.

## 6.4 LATENCY PERFORMANCE

In cloud environments, latency is an important consideration for anomaly detection. While CloudAnoAgent does exhibit higher detection latency than traditional ML baselines, this difference reflects a fundamental performance–latency trade-off. Classical ML models provide very low latency but suffer from high false positive rates, which in real deployments lead to operational overhead, alert fatigue, and inefficient use of engineering time. By contrast, CloudAnoAgent achieves

| Model           | Mean $\pm$ Std Dev    | Max      |
|-----------------|-----------------------|----------|
| IsolationForest | 0.61 ms $\pm$ 0.59 ms | 4.31 ms  |
| KMeans          | 0.25 ms $\pm$ 0.46 ms | 5.54 ms  |
| RarityModel     | 1.59 ms $\pm$ 0.78 ms | 12.76 ms |
| OOV Detector    | 1.15 ms $\pm$ 0.68 ms | 19.07 ms |

Table 5: Latency of ML Algorithms

| Model                 | Total Latency |
|-----------------------|---------------|
| Gemini 2.5 Flash      | 7.514 s       |
| Gemini 2.5 Flash Lite | 0.577 s       |
| Gemini 2.0 Flash      | 0.671 s       |
| GPT-4o                | 2.917 s       |
| GPT-4o Mini           | 1.834 s       |
| Qwen Plus             | 2.761 s       |

Table 6: Latency of Vanilla Prompting

| Model                 | Metrics Agent | Log Agent | Integrated Agent | Symbolic Verifier | Total   |
|-----------------------|---------------|-----------|------------------|-------------------|---------|
| Gemini 2.5 Flash      | 3.78 s        | 2.91 s    | 3.45 s           | 0.42 s            | 10.56 s |
| Gemini 2.5 Flash Lite | 0.63 s        | 0.59 s    | 0.71 s           | 0.42 s            | 2.35 s  |
| Gemini 2.0 Flash      | 1.08 s        | 0.82 s    | 1.05 s           | 0.39 s            | 3.34 s  |
| GPT-4o                | 2.33 s        | 1.92 s    | 2.09 s           | 0.41 s            | 6.75 s  |
| GPT-4o Mini           | 1.35 s        | 1.09 s    | 1.23 s           | 0.40 s            | 4.07 s  |
| Qwen Plus             | 5.67 s        | 2.75 s    | 5.29 s           | 0.42 s            | 14.14 s |

Table 7: Latency of CloudAnoAgent

substantially lower false positive rates and can identify the specific anomaly scenario, enabling faster and more accurate root-cause analysis. Additionally, end-to-end latency is not fixed; it varies significantly depending on factors such as API provider, model size, and whether the LLM is locally hosted. Using smaller locally deployed models can further reduce latency when needed.

## 7 CONCLUSION AND FUTURE WORK

In this work, we introduce CLOUDANOAGENT, a large-scale benchmark for context-level anomaly detection in cloud environments. Unlike existing datasets, our benchmark jointly incorporates both metrics and logs, and provides annotations not only for anomaly presence but also for specific anomaly scenarios. It covers 29 anomalous scenarios and further includes challenging deceptive normal cases designed to test the false positive robustness of detection methods. To address these challenges, we propose CLOUDANOAGENT, a symbolic-verification guided LLM-based multi-agent system capable of handling multimodal data for anomaly detection. Experimental results demonstrate that CLOUDANOAGENT significantly outperforms baselines in both anomaly detection and anomaly scenario identification, while also showing promising generalization to point anomaly datasets of different modalities. Future work will explore advanced post-training strategies and LLM architectural improvements to further enhance scenario identification accuracy, as well as optimizations to reduce detection latency in real-world deployments.

## REFERENCES

- Amira Mahamat Abdallah, Aysha Saif Rashed Obaid Alkaabi, Ghaya Bark Nasser Douman Alameri, Saida Hafsa Rafique, Nura Shifa Musa, and Thangavel Murugan. Cloud network anomaly detection using machine and deep learning techniques— recent research advancements. *IEEE Access*, 12:56749–56773, 2024. doi: 10.1109/ACCESS.2024.3390844.
- Lahiru Akmeemana, Chamodya Attanayake, Husni Faiz, and Sandareka Wickramanayake. Gal-mad: Towards explainable anomaly detection in microservice applications using graph attention networks. *arXiv preprint arXiv:2504.00058*, 2025.
- Alexander Bakhtin, Jesse Nyysölä, Yuqing Wang, Noman Ahmad, Ke Ping, Matteo Esposito, Mika Mäntylä, and Davide Taibi. Lo2: Microservice api anomaly dataset of logs and metrics. In *Proceedings of the 21st International Conference on Predictive Models and Data Analytics in Software Engineering*, pp. 1–10, 2025.
- Sakil Barbhuiya, Zafeirios Papazachos, Peter Kilpatrick, and Dimitrios S Nikolopoulos. Rads: Real-time anomaly detection system for cloud data centres. *arXiv preprint arXiv:1811.04481*, 2018a.

- 
- 540 Sakil Barbhuiya, Zafeirios C. Papazachos, Peter Kilpatrick, and Dimitrios S. Nikolopoulos. Rads:  
541 Real-time anomaly detection system for cloud data centres. *ArXiv*, abs/1811.04481, 2018b. URL  
542 <https://api.semanticscholar.org/CorpusID:53281451>.  
543
- 544 Peter Bodik, Armando Fox, Michael J. Franklin, Michael I. Jordan, and David A. Patterson. Char-  
545 acterizing, modeling, and generating workload spikes for stateful services. In *Proceedings of the*  
546 *1st ACM Symposium on Cloud Computing*, SoCC '10, pp. 241–252, New York, NY, USA, 2010.  
547 Association for Computing Machinery. ISBN 9781450300360. doi: 10.1145/1807128.1807166.  
548 URL <https://doi.org/10.1145/1807128.1807166>.
- 549 Jiahao Bu, Ying Liu, Shenglin Zhang, Weibin Meng, Qitong Liu, Xiaotian Zhu, and Dan Pei. Rapid  
550 deployment of anomaly detection models for large number of emerging kpi streams. In *2018 IEEE*  
551 *37th International Performance Computing and Communications Conference (IPCCC)*, pp. 1–8,  
552 2018. doi: 10.1109/PCCC.2018.8711315.
- 553 Peter Garraghan, Paul Townend, and Jie Xu. An empirical failure-analysis of a large-scale cloud  
554 computing environment. In *2014 IEEE 15th International Symposium on High-Assurance Systems*  
555 *Engineering*, pp. 113–120, 2014. doi: 10.1109/HASE.2014.24.
- 556 Wei Guan, Jian Cao, Shiyong Qian, Jianqi Gao, and Chun Ouyang. Logllm: Log-based anomaly  
557 detection using large language models. *arXiv preprint arXiv:2411.08561*, 2024.  
558
- 559 Haryadi S. Gunawi, Mingzhe Hao, Tanakorn Leesatapornwongsa, Tiratat Patana-anake, Thanh Do,  
560 Jeffry Adityatama, Kurnia J. Eliazar, Agung Laksono, Jeffrey F. Lukman, Vincentius Martin,  
561 and Anang D. Satria. What bugs live in the cloud? a study of 3000+ issues in cloud systems.  
562 In *Proceedings of the ACM Symposium on Cloud Computing*, SOCC '14, pp. 1–14, New York,  
563 NY, USA, 2014. Association for Computing Machinery. ISBN 9781450332521. doi: 10.1145/  
564 2670979.2670986. URL <https://doi.org/10.1145/2670979.2670986>.
- 565 Haixuan Guo, Shuhan Yuan, and Xintao Wu. Logbert: Log anomaly detection via bert. In *2021*  
566 *International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2021. doi: 10.1109/  
567 IJCNN52387.2021.9534113.
- 568 Kate Highnam, Kai Arulkumaran, Zachary Hanif, and Nicholas R Jennings. Beth dataset: Real cy-  
569 bersecurity data for unsupervised anomaly detection research. In *CEUR Workshop Proc*, volume  
570 3095, pp. 1–12, 2021.  
571
- 572 Mohammad S. Islam, William Pourmajidi, Lei Zhang, John Steinbacher, Tony Erwin, and Andriy  
573 Miranskyy. Anomaly detection in a large-scale cloud platform. In *2021 IEEE/ACM 43rd Interna-*  
574 *tional Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*, pp.  
575 150–159, 2021. doi: 10.1109/ICSE-SEIP52600.2021.00024.
- 576 Mohammad Saiful Islam, Mohamed Sami Rakha, William Pourmajidi, Janakan Sivaloganathan,  
577 John Steinbacher, and Andriy Miranskyy. Anomaly detection in large-scale cloud systems: An  
578 industry case and dataset. In *2025 IEEE/ACM 47th International Conference on Software Engi-*  
579 *neering: Software Engineering in Practice (ICSE-SEIP)*, pp. 377–388. IEEE, 2025.  
580
- 581 Reza Sherafat Kazemzadeh and Hans-Arno Jacobsen. Reliable and highly available distributed  
582 publish/subscribe service. In *2009 28th IEEE International Symposium on Reliable Distributed*  
583 *Systems*, pp. 41–50, 2009. doi: 10.1109/SRDS.2009.32.
- 584 Shalini Kumari, Chander Prabha, Asif Karim, Md. Mehedi Hassan, Sami Azam, and Peican Zhu. A  
585 comprehensive investigation of anomaly detection methods in deep learning and machine learn-  
586 ing: 2019–2023. *IET Information Security*, 2024, November 2024. doi: 10.1049/2024/8821891.  
587 URL <https://doi.org/10.1049/2024/8821891>.
- 588 Qingwei Lin, Hongyu Zhang, Jian-Guang Lou, Yu Zhang, and Xuwei Chen. Log clustering based  
589 problem identification for online service systems. In *2016 IEEE/ACM 38th International Confer-*  
590 *ence on Software Engineering Companion (ICSE-C)*, pp. 102–111, 2016.  
591
- 592 Ping Lou, Yun Yang, and Junwei Yan. An anomaly detection method for cloud service platform.  
593 In *Proceedings of the 2019 4th International Conference on Machine Learning Technologies*,  
ICMLT '19, pp. 70–75, New York, NY, USA, 2019. Association for Computing Machinery. ISBN

---

594 9781450363235. doi: 10.1145/3340997.3341005. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3340997.3341005)  
595 3340997.3341005.  
596

597 Adetokunbo A.O. Makanju, A. Nur Zincir-Heywood, and Evangelos E. Milios. Clustering event logs  
598 using iterative partitioning. In *Proceedings of the 15th ACM SIGKDD International Conference*  
599 *on Knowledge Discovery and Data Mining, KDD '09*, pp. 1255–1264, New York, NY, USA,  
600 2009. Association for Computing Machinery. ISBN 9781605584959. doi: 10.1145/1557019.  
601 1557154. URL <https://doi.org/10.1145/1557019.1557154>.

602 Weibin Meng, Ying Liu, Yichen Zhu, Shenglin Zhang, Dan Pei, Yuqing Liu, Yihao Chen, Ruizhi  
603 Zhang, Shimin Tao, Pei Sun, and Rong Zhou. Loganomaly: Unsupervised detection of sequential  
604 and quantitative anomalies in unstructured logs. In *Proceedings of the Twenty-Eighth Interna-*  
605 *tional Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 4739–4745. International Joint  
606 Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/658. URL  
607 <https://doi.org/10.24963/ijcai.2019/658>.

608 Animesh Nandi, Atri Mandal, Shubham Atreja, Gargi B. Dasgupta, and Subhrajit Bhattacharya.  
609 Anomaly detection using program control flow graph mining from execution logs. In *Proceed-*  
610 *ings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data*  
611 *Mining, KDD '16*, pp. 215–224, New York, NY, USA, 2016. Association for Computing Machin-  
612 ery. ISBN 9781450342322. doi: 10.1145/2939672.2939712. URL [https://doi.org/10.](https://doi.org/10.1145/2939672.2939712)  
613 [1145/2939672.2939712](https://doi.org/10.1145/2939672.2939712).

614 Adam Oliner and Jon Stearley. What supercomputers say: A study of five system logs. In *37th*  
615 *Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN'07)*,  
616 pp. 575–584, 2007. doi: 10.1109/DSN.2007.103.

617 Ron Sun and Lawrence A. Bookman (eds.). *Computational architectures integrating neural and*  
618 *symbolic processes: A perspective on the state of the art*. Springer International Series in En-  
619 gineering and Computer Science, vol. 260. Kluwer Academic Publishers / Springer, 1994. doi:  
620 10.1007/b102608.  
621

622 Wei Xu, Ling Huang, Armando Fox, David Patterson, and Michael I Jordan. Detecting large-scale  
623 system problems by mining console logs. In *Proceedings of the ACM SIGOPS 22nd symposium*  
624 *on Operating systems principles*, pp. 117–132, 2009.

625 Shenglin Zhang, Ying Liu, Dan Pei, Yu Chen, Xianping Qu, Shimin Tao, and Zhi Zang. Rapid and  
626 robust impact assessment of software changes in large internet-based services. In *Proceedings of*  
627 *the 11th ACM Conference on Emerging Networking Experiments and Technologies, CoNEXT '15*,  
628 New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450334129. doi:  
629 10.1145/2716281.2836087. URL <https://doi.org/10.1145/2716281.2836087>.  
630

631 Shenglin Zhang, Ying Liu, Dan Pei, Yu Chen, Xianping Qu, Shimin Tao, Zhi Zang, Xiaowei Jing,  
632 and Mei Feng. Funnel: Assessing software changes in web-based services. *IEEE Transactions*  
633 *on Services Computing*, 11(1):34–48, 2018. doi: 10.1109/TSC.2016.2539945.

634 Jingwen Zhou, Zhenbang Chen, Ji Wang, Zibin Zheng, and Michael R. Lyu. Tracebench: An open  
635 data set for trace-oriented monitoring. In *Proceedings of the 6th IEEE International Conference*  
636 *on Cloud Computing Technology and Science (CloudCom 2014)*, pp. 519–526. IEEE Computer  
637 Society, 2014. doi: 10.1109/CloudCom.2014.79. URL [http://mtracer.github.io/](http://mtracer.github.io/TraceBench/)  
638 [TraceBench/](http://mtracer.github.io/TraceBench/).  
639  
640  
641  
642  
643  
644  
645  
646  
647

## A LIMITATIONS

Because CLOUDANOAGENT incorporates LLMs in both anomaly detection and scenario reasoning, its performance is not entirely deterministic. Subtle differences in prompt wording or changes in model versions may lead to slightly different outputs, introducing variability across repeated experiments. Similarly, since CLOUDANOBENCH contains partially LLM-synthesized data, the dataset itself may be sensitive to model behavior and prompt design. While we carefully design prompts and incorporate symbolic checks to mitigate these issues, achieving strict reproducibility under evolving LLM behavior remains challenging.

## B ABLATION STUDY FOR METRICS AGENT ONLY AND LOG AGENT ONLY

| Model                 | P     | R     | FPR   | F1    | SI:Acc |
|-----------------------|-------|-------|-------|-------|--------|
| Gemini 2.5 Flash      | 58.55 | 93.27 | 79.26 | 71.94 | 0.426  |
| Gemini 2.5 Flash Lite | 59.31 | 93.70 | 77.15 | 72.64 | 0.407  |
| Gemini 2.0 Flash Lite | 57.57 | 87.99 | 77.86 | 69.60 | 0.353  |
| GPT-4o Mini           | 57.14 | 86.68 | 78.03 | 68.88 | 0.387  |
| GPT-4o                | 59.04 | 89.90 | 74.87 | 71.27 | 0.411  |
| Qwen Plus             | 59.10 | 90.34 | 75.04 | 71.45 | 0.382  |

Table 8: Ablation Study: CloudAnoAgent (Metrics Agent Only)

| Model                 | P     | R     | FPR   | F1    | SI:Acc |
|-----------------------|-------|-------|-------|-------|--------|
| Gemini 2.5 Flash      | 76.28 | 85.21 | 31.81 | 80.50 | 0.543  |
| Gemini 2.5 Flash Lite | 76.93 | 84.48 | 30.40 | 80.53 | 0.518  |
| Gemini 2.0 Flash Lite | 78.55 | 80.97 | 26.54 | 79.74 | 0.452  |
| GPT-4o Mini           | 78.68 | 80.53 | 26.19 | 79.59 | 0.539  |
| GPT-4o                | 78.17 | 83.89 | 28.12 | 80.93 | 0.574  |
| Qwen Plus             | 78.13 | 82.14 | 27.59 | 80.09 | 0.553  |

Table 9: Ablation Study: CloudAnoAgent (Log Agent Only)

## C ALIGNMENT BETWEEN CLOUDANOBENCH AND REAL-WORLD DISTRIBUTION

To examine whether CLOUDANOBENCH reflects realistic temporal behaviors rather than synthetic artifacts, we performed an alignment study comparing our benchmark with the RS-Anomic dataset (Akmeemana et al., 2025). RS-Anomic contains multivariate system metrics collected from the RobotShop Instana microservice application and is widely used for anomaly detection and root-cause localization in production-like environments.

Since direct comparison of logs is not meaningful due to system-specific identifiers and environment-dependent log formats, our analysis focuses on temporal similarity between system metrics. We selected two anomaly categories in CLOUDANOBENCH (CPU and networking) and matched them with their closest counterparts in RS-Anomic.

### C.1 SIMILARITY METRICS

We applied three complementary sequence-similarity metrics:

- **Spearman’s  $\rho$** : measures similarity in monotonic trends.
- **Kendall’s  $\tau$** : evaluates consistency in pairwise rank ordering.
- **DTW/len**: normalized Dynamic Time Warping distance, which quantifies the temporal deformation needed to align two sequences.

These metrics jointly capture trend agreement, local ordering consistency, and temporal-shape similarity.

| Category | Scenario                 | Spearman $\rho$ | Kendall $\tau$ | DTW/len      |
|----------|--------------------------|-----------------|----------------|--------------|
| CPU      | High CPU Usage / CPU Hog | <b>0.816</b>    | <b>0.656</b>   | <b>0.179</b> |
| Network  | High Latency             | <b>0.688</b>    | <b>0.513</b>   | <b>0.169</b> |
| Network  | Out-of-order Packets     | <b>0.659</b>    | <b>0.496</b>   | <b>0.160</b> |
| Network  | Packet Loss              | <b>0.659</b>    | <b>0.494</b>   | <b>0.182</b> |

Table 10: Similarity between CLOUDANOBENCH and RS-Anomic metric traces across matched anomaly scenarios.

## C.2 FINDINGS

Across CPU and multiple network anomalies, we observe consistently strong alignment:

- **CPU anomalies** show high agreement in global trends and temporal patterns ( $\rho = 0.816$ ,  $\tau = 0.656$ ,  $\text{DTW}/\text{len} = 0.179$ ).
- **Network anomalies** also demonstrate close correspondence to real-world dynamics, with  $\rho \approx 0.66$ – $0.69$  and  $\text{DTW}/\text{len} \approx 0.16$ – $0.18$ .

These quantitative results are supported by qualitative inspection. For matched scenarios, both datasets show characteristic rise–fall dynamics, temporal jitter patterns, and noise structure. CPU hog events display rapid saturation and sustained utilization, while network-related anomalies exhibit latency buildup, irregular spikes, and transient bursts.

Although CLOUDANOBENCH is constructed using LLM-assisted scenario design and code-driven data generation, the resulting metric traces follow real-world temporal signatures closely. The strong correspondence observed across both quantitative metrics and qualitative patterns supports the realism of CLOUDANOBENCH and reinforces the validity of evaluating CloudAnoAgent on this benchmark.

| Scenario                                 | Metrics  | Logs   |
|--|--|--|
| CPU Hog Process                          | CPU spike  | Timeout errors   |
| Memory Leak                              | CPU fluctuation, Memory increase, Disk I/O spike               | OOM killer   |
| Disk I/O Bottleneck                      | CPU iowait spike, Disk I/O saturation, Network dip/fluctuation | Task blocked, Slow queries   |
| Noisy Neighbor                           | CPU steal spike  | Unexplained slowdowns, Timeouts  |
| Handle/Thread Exhaustion                 | CPU & Memory spike, Network dip                                | Resource allocation errors   |
| High Network Latency / Packet Loss       | CPU spike, Memory increase, Network dip                        | Connection errors  |
| Bandwidth Saturation                     | CPU spike, Network cap plateau                                 | Timeout errors, Packet loss  |
| DNS Resolution Failure                   | Network drop   | DNS resolution errors  |
| Firewall/Security Group Misconfiguration | Network drop (blocked ports)                                   | Blocked port timeouts  |
| Application Crash Loop                   | CPU & Memory fluctuation, I/O fluctuation                      | Service crash/restart, Fatal errors  |
| Deadlock / Livelock                      | CPU dip/spike, Throughput drop                                 | Deadlock detected  |
| External Dependency Failure              | CPU dip/spike, Memory increase, Network dip                    | Dependency failures, 503 errors  |
| Application Misconfiguration             | CPU & Memory spike, Network dip, Fluctuations                  | Misconfiguration errors  |
| Garbage Collection (GC) Storm            | CPU spike, Memory sawtooth fluctuation, Throughput dip         | GC pauses, Timeouts  |
| Cryptomining Malware                     | CPU & GPU spike, Network fluctuation                           | No matching logs, Unknown process  |
| DDoS Botnet Agent                        | CPU & Memory spike, Network saturation                         | Resource contention, Outbound flood  |
| Ransomware                               | CPU & Memory spike, Disk I/O spike                             | File errors, Log clearing, Ransom note                                       |
| Data Exfiltration                        | CPU & Memory & Disk I/O & Network spike                        | Abnormal queries, File reads, Data exfiltration                              |
| Spam Bot                                 | CPU & Network spike  | Spam activity, Bounce messages   |
| Rootkit / Backdoor                       | N/A  | Log tampering, Rootkit traces  |
| Brute-force / Credential Stuffing Attack | CPU & Network spike  | Authentication failures, Log storm   |
| Web Shell Beaconing & Command Execution  | CPU & Disk I/O spike   | Suspicious repeated requests; unexpected process creation                    |
| Password Cracking                        | CPU & GPU spike, Memory spike, Disk I/O startup spike          | Cracking process detected; bash/auth log evidence                            |
| Log Deletion / Tampering                 | CPU spike, Disk write spike                                    | Logs removed/cleared; file-integrity alerts; possible bash history           |
| Reverse Shell                            | CPU spikes, Persistent outbound connection                     | Unexpected outbound process; nonstandard outbound TCP in firewall/audit logs |
| Application-Layer DDoS Attack            | CPU & Memory & Disk I/O spike                                  | Endpoint request flood; DB slow queries and timeouts                         |
| Time Skew (Clock Drift)                  | N/A  | Clock skew errors  |
| Kernel / Driver Bug                      | Unpredictable spikes/fluctuations or drop to zero              | Kernel panic, Segfaults, HW faults   |

Table 11: Anomaly scenarios and their performance in metrics and logs.

## D ANOMALY SCENARIOS DETAILS

| Anomaly Scenario                         | Description  |
|--|--|
| CPU Hog Process                          | A single or few processes continuously occupy extremely high CPU due to bugs or malicious behavior         |
| Memory Leak                              | The application fails to release memory due to bugs, causing continuous memory depletion                   |
| Disk I/O Bottleneck                      | Disk read and write requests exceed hardware capacity  |
| Noisy Neighbor                           | Other VMs on the same physical host overconsume shared resources in a virtualized environment              |
| Handle/Thread Exhaustion                 | Processes open too many file handles or threads, reaching the system limit                                 |
| High Network Latency / Packet Loss       | Network congestion or failures cause packet delay or loss  |
| Bandwidth Saturation                     | Outbound or inbound traffic reaches the cloud provider's bandwidth cap                                     |
| DNS Resolution Failure                   | Domain names cannot be resolved to IP addresses  |
| Firewall/Security Group Misconfiguration | Incorrect firewall or cloud security group rules block normal port connections                             |
| Application Crash Loop                   | Applications repeatedly crash due to fatal errors and are restarted by a watchdog                          |
| Deadlock / Livelock                      | Two or more threads/processes wait for each other's resources, halting progress (deadlock/livelock)        |
| External Dependency Failure              | Dependent databases, caches, or third-party APIs become unavailable  |
| Application Misconfiguration             | An incorrect configuration file causes app startup failures, functional errors, or performance degradation |
| Garbage Collection (GC) Storm            | Excessive or long garbage collection pauses in managed-language apps (e.g., Java)                          |
| Cryptomining Malware                     | Attacker runs malicious program on the server to mine cryptocurrency (uses CPU/GPU).                       |
| DDoS Botnet Agent                        | Server used as a "bot" in a DDoS campaign, sending large volumes of attack traffic.                        |
| Ransomware                               | Ransomware encrypts files on the server and demands ransom.  |
| Data Exfiltration                        | Attacker exfiltrates sensitive data from the server to an external host.                                   |
| Spam Bot                                 | Server used to send large volumes of spam or phishing emails.  |
| Rootkit / Backdoor                       | Attacker installs stealthy persistent backdoor/rootkit to maintain covert access.                          |
| Brute-force / Credential Stuffing Attack | Ongoing automated credential-stuffing / brute-force attack against authentication services (SSH/RDP/Web).  |
| Web Shell Beaconing & Command Execution  | Web shell on server beaconing to C2 and executing commands   |
| Password Cracking                        | Offline password cracking with john/hashcat  |
| Log Deletion / Tampering                 | Log cleaning to cover tracks   |
| Reverse Shell                            | Reverse shell established to attacker host   |
| Application-Layer DDoS Attack            | Application-level DoS via expensive functions  |
| Time Skew (Clock Drift)                  | Server system time drifts from NTP standard  |
| Kernel / Driver Bug                      | OS kernel or hardware driver bug is triggered, causing instability or crashes                              |

Table 12: Anomaly scenarios with descriptions

## E DECEPTIVE NORMAL SCENARIOS DETAILS

| Normal Scenario                  | Description  |
|----------------------------------|--|
| Nightly Full Backup              | Nightly full backup with database dump and remote upload |
| Log Rotation & Compression       | Scheduled log rotation and compression                   |
| Periodic Data Aggregation        | Periodic data aggregation for reports                    |
| Filesystem Cleanup               | Filesystem cleanup of expired files                      |
| SSL Certificate Renewal          | Automatic SSL certificate renewal                        |
| Blue-Green Deployment            | Blue-green deployment traffic switch                     |
| Live Streaming Event Start       | Start of a large-scale live streaming event              |
| Search Engine Aggressive Crawl   | Aggressive crawl by search engine bots                   |
| End-of-Month Financial Closing   | End-of-month financial closing workload                  |
| On-Demand Full Report Generation | On-demand full report generation                         |
| ETL Pipeline Execution           | ETL pipeline execution with data extraction and loading  |
| On-Demand Video Transcoding      | On-demand video transcoding with FFmpeg                  |
| Database Compaction/Vacuum       | Database compaction or vacuuming                         |
| Cache Eviction Storm             | Cache eviction storm with many misses                    |
| Connection Pool Scaling          | Dynamic expansion of DB connection pool                  |
| File Integrity Monitoring Scan   | File integrity monitoring scan                           |

Table 14: Normal scenarios with descriptions.

| Normal Scenario                  | Metrics  | Logs                                 |
|----------------------------------|--|--------------------------------------|
| Nightly Full Backup              | CPU spike (compression), Disk I/O spike (read), Net out spike (upload) | Backup task logs                     |
| Log Rotation & Compression       | CPU spike (compression), Disk I/O spike                                | Logrotate success logs               |
| Periodic Data Aggregation        | CPU spike (scheduled), Disk I/O spike                                  | Report job logs, DB slow query logs  |
| Filesystem Cleanup               | CPU spike (scan), Disk I/O spike (read/write)                          | Cleanup script logs                  |
| SSL Certificate Renewal          | Net in/out spike (CA communication), CPU spike (reload)                | Certbot/ACME renewal logs            |
| Blue-Green Deployment            | Green env: CPU/Mem/Net spike; Blue env: CPU/Mem/Net dip                | Load balancer & deployment logs      |
| Live Streaming Event Start       | Net out spike, Connection spike, CPU spike                             | Live stream logs, CDN access logs    |
| Search Engine Aggressive Crawl   | Net in/out spike, CPU spike, Disk I/O spike                            | Web access logs (bot user-agent)     |
| End-of-Month Financial Closing   | CPU/Mem/Disk I/O gradual increase & fluctuation                        | App logs with financial keywords     |
| On-Demand Full Report Generation | CPU spike, Memory increase, Disk I/O spike                             | Report generation logs               |
| ETL Pipeline Execution           | CPU/Mem/Disk/Net spikes & fluctuations                                 | ETL tool logs (Airflow/dbt etc.)     |
| On-Demand Video Transcoding      | Sustained CPU spike, Disk I/O spike                                    | Transcoding job logs (FFmpeg)        |
| Database Compaction/Vacuum       | Disk I/O spike, CPU spike  | DB compaction/vacuum logs            |
| Cache Eviction Storm             | DB CPU/Disk I/O/Net in spike, Cache misses spike                       | Cache miss logs, DB query surge logs |
| Connection Pool Scaling          | App latency spike, CPU spike   | Connection pool expansion logs       |
| File Integrity Monitoring Scan   | Disk I/O spike (read), CPU spike (hashing)                             | AIDE/Tripwire scan logs              |

Table 13: Normal scenarios and their performance in metrics and logs.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

F BENCHMARK EXAMPLES

| timestamp            | cpu_usage | mem_usage | disk_io | net_in | net_out |
|----------------------|-----------|-----------|---------|--------|---------|
| 2025-07-12T04:00:00Z | 13.75     | 36.2      | 20.12   | 1.31   | 0.89    |
| 2025-07-12T04:00:05Z | 19.51     | 42.13     | 16.7    | 1.31   | 1.14    |
| 2025-07-12T04:00:10Z | 17.32     | 42.61     | 28.87   | 1.37   | 0.96    |
| 2025-07-12T04:00:15Z | 15.99     | 40.61     | 28.16   | 1.41   | 1.05    |
| 2025-07-12T04:00:20Z | 11.56     | 42.71     | 18.87   | 1.01   | 1.44    |
| 2025-07-12T04:00:25Z | 11.56     | 39.94     | 24.9    | 1      | 0.89    |
| 2025-07-12T04:00:30Z | 10.58     | 40.23     | 27.26   | 1.3    | 1.46    |
| 2025-07-12T04:00:35Z | 18.66     | 39.28     | 23.33   | 1.15   | 1.41    |
| 2025-07-12T04:00:40Z | 16.01     | 35.25     | 22.94   | 1.2    | 0.7     |
| 2025-07-12T04:00:45Z | 17.08     | 36.08     | 18.63   | 1.3    | 0.57    |
| 2025-07-12T04:00:50Z | 10.21     | 35.31     | 16.4    | 1.39   | 0.6     |
| 2025-07-12T04:00:55Z | 19.7      | 41.36     | 28.46   | 0.84   | 0.52    |
| 2025-07-12T04:01:00Z | 18.32     | 38.14     | 28.51   | 0.88   | 0.59    |
| 2025-07-12T04:01:05Z | 12.12     | 40.09     | 24.5    | 0.59   | 1.18    |
| 2025-07-12T04:01:10Z | 11.82     | 44.08     | 20.09   | 1.08   | 0.57    |
| 2025-07-12T04:01:15Z | 11.83     | 37.49     | 20.24   | 0.54   | 0.82    |
| 2025-07-12T04:01:20Z | 13.04     | 39.1      | 25.89   | 0.97   | 1.34    |
| 2025-07-12T04:01:25Z | 15.25     | 42.56     | 28.46   | 1.04   | 0.52    |
| 2025-07-12T04:01:30Z | 14.32     | 37.29     | 28.31   | 0.79   | 1.31    |
| 2025-07-12T04:01:35Z | 12.91     | 35.77     | 26.7    | 1.09   | 0.78    |
| 2025-07-12T04:01:40Z | 16.12     | 37.9      | 24.63   | 0.53   | 0.62    |
| 2025-07-12T04:01:45Z | 11.39     | 36.61     | 16.26   | 0.54   | 1.2     |
| 2025-07-12T04:01:50Z | 12.92     | 44.3      | 17.42   | 1.32   | 1.13    |
| 2025-07-12T04:01:55Z | 13.66     | 43.08     | 28.48   | 0.86   | 1.38    |
| 2025-07-12T04:02:00Z | 14.56     | 41.33     | 24.1    | 0.63   | 1.24    |
| 2025-07-12T04:02:05Z | 17.85     | 43.71     | 15.14   | 1.02   | 1.3     |
| 2025-07-12T04:02:10Z | 12        | 43.04     | 16.52   | 1.27   | 0.78    |
| 2025-07-12T04:02:15Z | 15.14     | 36.87     | 24.95   | 0.72   | 0.68    |
| 2025-07-12T04:02:20Z | 15.92     | 43.93     | 15.08   | 1.12   | 1.25    |
| 2025-07-12T04:02:25Z | 10.46     | 40.39     | 17.41   | 0.59   | 1.31    |

Figure 4: Metrics Example of Data Exfiltration (Simplified Version)

```

972 Jul 12 04:00:10 systemd[1]: Started Session 5 of user db_admin.
973 Jul 12 04:00:25 sshd[15482]: Accepted publickey for db_admin from 10.1.2.5 port 49822 ssh2
974 Jul 12 04:00:50 cron[11432]: (root) CMD (run-parts --report /etc/cron.hourly)
975 Jul 12 04:01:40 web-app[887]: INFO: User 'alex' successfully authenticated.
976 Jul 12 04:02:20 web-app[887]: GET /api/v1/dashboard HTTP/1.1 200
977 Jul 12 04:02:35 auditd[512]: type=SYSCALL msg=audit(...): arch=c000003e syscall=257 success=yes
978 exit=3 a0=ffffff9c a1=7fffd1f9b6e20 a2=90800 a3=0 items=1 ppid=15510 pid=15512 auid=1001 uid=1001
979 gid=1001 euid=1001 suid=1001 fsuid=1001 egid=1001 sgid=1001 fsgid=1001 tty=pts0 ses=5 comm="cat"
980 exe="/bin/cat" key="secret_file_access"
981 Jul 12 04:02:40 auditd[512]: type=PATH msg=audit(...): item=0 name="/etc/shadow" inode=131078
982 dev=fd:01 mode=0100640 ouid=0 ogid=42 rdev=00:00 nametype=NORMAL
983 Jul 12 04:02:50 auditd[512]: type=SYSCALL msg=audit(...): arch=c000003e syscall=257 success=yes
984 exit=4 a0=ffffff9c a1=7fffd1f9b6e20 a2=90800 a3=0 items=1 ppid=15510 pid=15515 auid=1001 uid=1001
985 gid=1001 euid=1001 suid=1001 fsuid=1001 egid=1001 sgid=1001 fsgid=1001 tty=pts0 ses=5 comm="cat"
986 exe="/bin/cat" key="secret_file_access"
987 Jul 12 04:02:55 auditd[512]: type=PATH msg=audit(...): item=0 name="/opt/app/credentials.env"
988 inode=262148 dev=fd:01 mode=0100600 ouid=1001 ogid=1001 rdev=00:00 nametype=NORMAL
989 Jul 12 04:03:10 mysqld[1234]: 2025-07-12T04:03:10.123456Z 115 [Note] User 'db_admin'@'localhost'
990 executed: 'SELECT * FROM customers.credit_card_details;'
991 Jul 12 04:03:30 mysqld[1234]: 2025-07-12T04:03:30.567890Z 115 [Note] Query OK, 1572880 rows in set
992 (28.15 sec)
993 Jul 12 04:03:50 ufw-audit[15601]: [UFW ALLOW] IN= OUT=eth0 SRC=172.17.0.5 DST=185.125.190.44
994 LEN=60 TOS=0x00 PREC=0x00 TTL=64 ID=12345 DF PROTO=TCP SPT=48128 DPT=443 WINDOW=65535 RES=0x00 SYN
995 Jul 12 04:06:30 sshd[15482]: pam_unix(sshd:session): session closed for user db_admin
996 Jul 12 04:06:50 systemd[1]: Starting Clean up of User Sessions...
997 Jul 12 04:07:05 systemd[1]: Removed slice User Slice of UID 1001.
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

```

Figure 5: Log Example of Data Exfiltration (Simplified Version)

## G COMPARED DATASETS

**Blue Gene/L Dataset.** The Blue Gene/L dataset is provided by Lawrence Livermore National Laboratory (LLNL). The dataset comprises approximately 348,460 alert messages, spanning 41 distinct alert categories. The dataset covers a time period of 215 days, totaling approximately 1.207 GB of log data, and includes a wide range of system anomalies, from hardware failures to software errors.

**Thunderbird Dataset.** The Thunderbird dataset is provided by Sandia National Laboratories (SNL). The dataset comprises approximately 3,248,239 alert messages, spanning 10 distinct alert categories. The dataset covers a period of 244 days, totaling approximately 27.367 GB of log data, which includes a range of system anomalies such as hardware failures, operating system issues, and network and disk system errors.

**HPC Dataset.** The HPC dataset was collected from the Los Alamos National Laboratory, which comprises large-scale event logs from a high-performance computing cluster. Specifically, the dataset consists of 433,490 log messages and encompasses 106 distinct log formats. These logs capture various state transitions within the cluster's operations, including system configurations, environmental changes, and error messages. This dataset serves as a crucial benchmark for evaluating the performance of log analysis algorithms.

**HDFS Dataset.** The HDFS dataset was collected from a Hadoop cluster deployed on over 200 EC2 nodes, processing more than 200 TB of random data over a 48-hour period. The dataset consists of over 24 million log entries, capturing a wide range of system operations related to distributed storage and large-scale data processing. These logs include critical information on block allocation, replication, deletion, and other key operations, providing a valuable resource for performance monitoring and anomaly detection in large-scale distributed computing environments.

**BETH Dataset.** The BETH dataset was collected through the BPF-Extended Tracking Honeypot, containing over 8 million data records that encompass host activities and attack behaviors within a cloud environment. Specifically designed for cybersecurity anomaly detection research, the BETH dataset is primarily used for unsupervised anomaly detection tasks. The dataset includes kernel-level process invocation logs and network-level DNS query logs, with each record manually annotated to distinguish between benign and malicious behaviors. Logs for each host contain benign activity and

---

at most one attack event. The dataset is divided into training, validation, and test sets, making it suitable for performance evaluation of behavioral analysis and anomaly detection algorithms.

**IBM Cloud Console Dataset.** The IBM Cloud Console dataset is a large-scale telemetry dataset collected from the IBM Cloud Console over a period of 4.5 months, from January 22, 2024, to June 7, 2024. The dataset contains 39,365 rows and 117,448 columns, with each row representing a 5-minute time interval and capturing response time information from microservices across IBM Cloud data centers. Column names in the dataset represent various features such as timestamp, location, type, host ID, method, status code, and API endpoint. This dataset is intended for anomaly detection research, providing real-world cloud environment data that aids in the development and validation of anomaly detection methods applied to large-scale cloud computing systems.

**RS-Anomic Dataset.** The RS-Anomic dataset is a multivariate time-series dataset created based on the RobotShop microservice application, designed to support anomaly detection research in microservice architectures. The dataset comprises 100,464 normal instances and 14,112 anomalous instances, covering ten distinct types of anomalous behaviors, including service downtime, high concurrent user load, high CPU usage, packet loss, and response delay.

## H BASELINES

**LogBERT.** This work adopts a self-supervised Transformer that learns from normal-only log sequences via (i) masked log-key prediction (MLKP) and (ii) a hypersphere compactness loss (VHM) to encourage tight clusters of normal sequences. Raw logs are parsed into templates (e.g., *Drain*), then sequenced by dataset convention (*HDFS\_v1*: session-based; *BGL*: 5-minute sliding windows; *ThunderBird*: 1-minute windows). Following the public setup, the authors configure a 2-layer Transformer encoder (50-d input representations, 256-d hidden states). During inference, masked keys are predicted and an anomaly is declared when observed keys fall outside the top- $g$  candidates or when the count of unexpected keys exceeds a threshold  $r$ . This configuration yields strong precision/recall trade-offs across all three datasets, with particularly strong results on longer sequences in *ThunderBird*.

**LogAnomaly.** The authors employ an unsupervised LSTM-based detector that jointly captures (a) *sequential* patterns (via next-template prediction) and (b) *quantitative* invariants (via template count vectors). Templates are produced by a log parser (e.g., FT-Tree) and embedded with Template2Vec to inject semantic similarity (synonyms/antonyms). The default configuration uses a window size of 20 (step 1) and two LSTM layers (128 units each), trained on the earlier portion of logs (normal-only) and evaluated on the holdout. Modeling both sequence order and template-count relations reduces false alarms and achieves high F1 on *HDFS\_v1* and *BGL*; the combined view is especially beneficial when sequential cues alone are ambiguous.

**LogLLM.** This work treats windowed logs as natural-language context and applies an instruction/few-shot prompting scheme to (i) summarize state, (ii) assess deviations, and (iii) output an anomaly decision with rationale. The approach requires no dataset-specific fine-tuning; performance depends primarily on prompt design (structure, exemplars), windowing, and decoding settings (e.g., temperature). In the experiments on *HDFS\_v1*, *ThunderBird*, and *BGL*, the method attains competitive recall, while precision improves when template structure, light rules, or a symbolic post-verifier are incorporated to constrain spurious triggers.

**Isolation Forest.** Isolation Forest is a machine learning algorithm used for anomaly detection, which identifies outliers by constructing a set of random trees, known as "isolation trees." The algorithm works by randomly selecting features and split values to iteratively partition the data into smaller regions. Anomalies are generally easier to isolate, leading to shorter "isolation paths." Isolation Forest is particularly effective for high-dimensional data and offers high computational efficiency.

**Decision tree.** Decision tree is a commonly used machine learning algorithm for classification and regression tasks. It makes decisions by recursively partitioning the data into different regions. At

---

1080 each split, the data is divided into two subsets based on specific features and conditions, continuing  
1081 until a predetermined stopping criterion is met. The structure of a decision tree resembles a tree  
1082 diagram, starting from the root node and progressing through a series of splits, ultimately providing  
1083 a prediction at the leaf nodes. While decision trees are easy to understand and interpret, they are  
1084 prone to overfitting.

1085  
1086 **Logistic regression.** Logistic regression is a machine learning algorithm commonly used for bi-  
1087 nary classification problems. It computes the weighted sum of input features and applies a Sigmoid  
1088 function to convert the result into a probability value. This probability represents the likelihood of a  
1089 particular class, typically used to determine the class to which an instance belongs. Logistic regres-  
1090 sion is simple, easy to implement, and can output probability values, which are useful for assessing  
1091 the uncertainty of predictions.

1092  
1093 **K-Means.** K-Means is a widely used unsupervised learning algorithm for clustering analysis. It  
1094 partitions the data into several clusters, each represented by a centroid. The algorithm iteratively  
1095 adjusts the position of the centroids to ensure that data points within each cluster are as close as  
1096 possible to the centroid. This process continues until the cluster assignments no longer change.  
1097 K-Means is well-suited for large-scale data and offers high computational efficiency.

1098  
1099 **Rarity Model.** The Rarity Model is an unsupervised learning method based on data rarity, pri-  
1100 marily used for anomaly detection. The model identifies samples that occur less frequently in the  
1101 dataset and treats them as outliers. Typically, rare samples have a higher likelihood of being anoma-  
1102 lous compared to more common samples, making the Rarity Model effective in detecting these  
1103 anomalous data points.

1104  
1105 **OOV Detector.** The OOV Detector is a method used to identify "out-of-vocabulary" (OOV)  
1106 words. It works by recognizing words that were not present in the training data and marking them  
1107 as OOV. The OOV Detector is commonly applied in natural language processing tasks to handle  
1108 unseen vocabulary and enhance the model's ability to adapt to new words.

## 1109 I PROMPTS

```
1110  
1111  
1112  
1113 ### Anomalous Scenario: *1. CPU Hog Process*  
1114  
1115 **Description:**  
1116 One or a few processes continuously consume extremely high CPU resources due to a bug or malicious  
1117 behavior.  
1118  
1119 **Anomalous Metric Behavior:**  
1120 - **cpu usage:** `gradual increase` (sustained growth) or `spike` (sudden peak) up to 100% and  
1121 remains at that level.  
1122 - **gpu usage:** N/A  
1123 - **mem usage:** generally stable.  
1124 - **disk_io / net_in / net_out:** generally stable.  
1125  
1126 **Anomalous Log Behavior:**  
1127 - **Application logs:** request latency increases, with a large number of `timeout` errors.  
1128 - **System logs:** usually no directly related entries.  
1129  
1130 **Format Requirements:**  
1131 - CSV data must contain 90 rows (each row at 5-second intervals).  
1132 - Log data must be aligned with the CSV (timestamps must correspond).  
1133 - Log entries must use realistic formats and reflect real-world conditions.  
1134 - Include a small number of normal behavior log entries.  
1135 - Log files must not contain any direct descriptions of resource metrics.  
1136 - Both CSV and log files must resemble realistic operational data.  
1137 - File naming convention: `.csv` and `.log`.  
1138  
1139 **Note:** Two sample files are provided as references; please follow their format closely.
```

Figure 6: Prompt Example for Data Generation of CPU Hog Process (Simplified version)

---

```
1134 You need to decide whether it is "anomaly" or "normal". Set the "is_anomaly" is true if the result
1135 is anomaly, otherwise is false. Moreover, give the anomaly scenario, if it is normal choose "none".
1136 Please only output the answer **in following json format** without reason:
1137 ```json
1138 {
1139   "is_anomaly": true or false,
1140   "anomaly_scenario": a series of anomaly scenarios | "none"
1141 }
1142
```

Figure 7: Prompt Example for Vanilla LLM Prompting (Simplified version)

```
1143
1144
1145
1146
1147
1148
1149 You are a metrics anomaly detection agent.
1150
1151 Given a time series of metric data (e.g., CPU usage, memory usage, network traffic), determine:
1152 1. Whether the sequence contains an anomaly.
1153 2. If so, describe the anomaly and include its **approximate time of occurrence**.
1154 3. Identify which of the following five categories best describes the anomaly:
1155
1156   [1] Spike (sudden sharp increase)
1157   [2] Dip (sudden sharp decrease)
1158   [3] Gradual Increase (slow and steady rise)
1159   [4] Gradual Decrease (slow and steady fall)
1160   [5] Fluctuation (repeated high-variance oscillation)
1161
1162 Respond **only** with a JSON object in the following format:
1163
1164 ```json
1165 {
1166   "is_anomaly": true or false,
1167   "description": "A concise description of the anomaly, including approximate time, anomaly
1168 metrics and anomaly type in one sentence"
1169 }
1170 ```
1171
```

Figure 8: Prompt Example for Metrics Agent (Simplified version)

```
1172
1173
1174 You are a log anomaly detection agent.
1175
1176 Given a sequence of raw system or application log entries, your task is to:
1177 1. Analyze whether the logs indicate any abnormal or suspicious activity.
1178 2. Summarize the **key observed behavior in a single sentence**, and include the **time or time
1179 range** when the behavior occurred.
1180 3. Estimate the severity of the anomaly using the following levels:
1181   - "high": Indicators of intrusion, privilege abuse, data exfiltration, malware activity, or
1182   persistent service failures
1183   - "medium": Repeated login failures, unauthorized access, suspicious configuration changes
1184   - "low": Slight irregularities or uncommon but benign events
1185
1186 Respond **only** with a JSON object in the following format:
1187
1188 ```json
1189 {
1190   "description": "Summary of log behavior including time and anomaly type with anomaly event",
1191   "anomaly_level": "high" or "medium" or "low"
1192 }
1193 ```
1194
```

Figure 9: Prompt Example for Log Agent (Simplified version)

---

```

1188 You are an Integrated Agent responsible for evaluating whether a cloud system case represents
1189 a real anomaly by correlating metric and log data.
1190 You will be given:
1191 - The output of a metrics agent, in the form of:
1192   - `is_anomaly`: true or false
1193   - `description`: A short explanation including time and anomaly pattern
1194 - The output of a log agent, in the form of:
1195   - `anomaly_level`: "low", "medium", or "high"
1196   - `description`: A short explanation including time and observed behavior
1197 ---
1198 ### Your Task:
1199 1. Determine whether the system behavior represents a true anomaly.
1200 2. Use causal reasoning to explain how the metric pattern and log behavior relate to each
1201 other in time and implication.
1202 3. Provide a final conclusion with the following structure:
1203   - A boolean `is_anomaly` value.
1204   - If `is_anomaly` is true, classify the event by selecting one `anomaly_category` and one
1205     corresponding `anomaly_subtype` from the list below.
1206   - If `is_anomaly` is false, set both `anomaly_category` and `anomaly_subtype` to `none`.
1207   - A single-sentence `reason` that concisely combines the evidence from metrics and logs with
1208     timestamps and your logical interpretation.
1209 ---
1210 ### Anomaly Classification: a series of anomaly scenarios
1211 ---
1212 ### Special Attention:
1213 If the metrics agent reports an anomaly (`is_anomaly: true`) but the log agent reports
1214 `anomaly_level: "low"`:
1215 - Do not assume an anomaly by default.
1216 - First, check if the log entry provides a reasonable explanation for the metric behavior
1217 (e.g., a scheduled backup, a database indexing job, or a planned deployment).
1218 - If the log provides a benign explanation, treat the behavior as normal (`is_anomaly: false`).
1219 - Only report a true anomaly if the logs fail to plausibly explain the metric spike or if they
1220 contain subtle signs of early-stage risk despite the "low" level.
1221 Be conservative and prioritize avoiding false positives.
1222 ---
1223 ### Output Format (JSON only):
1224 Example for a anomaly case
1225 ```json
1226 {
1227   "is_anomaly": true,
1228   "anomaly_category": "resource_exhaustion_and_bottleneck",
1229   "anomaly_subtype": "memory_leak",
1230   "reason": "One-sentence causal explanation combining metrics and logs results."
1231 }
1232 Example for a non-anomaly:
1233 ```json
1234 {
1235   "is_anomaly": false,
1236   "anomaly_category": "none",
1237   "anomaly_subtype": "none",
1238   "reason": "The observed CPU spike at 14:32 UTC was caused by a scheduled daily database indexing
1239 job, which is considered normal system behavior."
1240 }
1241

```

Figure 10: Prompt Example for Integrated Agent (Simplified version)

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

## J SYMBOLIC VERIFICATION EXAMPLE FOR THE MINING SCENARIO (SIMPLIFIED)

Let the last 30 time steps be  $W = \{n - 29, \dots, n\}$ . Let the CPU threshold be  $\tau_{\text{cpu}} = 85$  and the minimum exceedance count be  $k_{\text{min}} = 24$ . Let the keyword set for mining be  $\mathcal{K}$  (e.g., xmrig, stratum+tcp, nanopool.org, custom-miner.service, update-run.sh, downloading payload, Accepted publickey, ufw). Let the minimum number of distinct keyword matches be  $s_{\text{min}} = 6$ . Denote the multivariate metric sequence by  $M$  and the multiset of log lines by  $L$ .

We write  $k \preceq_{\text{ci}} \ell$  to mean that keyword  $k \in \mathcal{K}$  appears as a (case-insensitive) substring of log line  $\ell \in L$ .

**Metric rule (sustained high CPU).** Define

$$C_{\text{cpu}}(M) = \sum_{i \in W} \mathbf{1}[M_i^{\text{cpu}} > \tau_{\text{cpu}}], \quad \text{VerifyMetric}_{\text{mining}}(M) = \mathbf{1}[C_{\text{cpu}}(M) \geq k_{\text{min}}]. \quad (3)$$

That is, the metric-side condition passes iff at least  $k_{\text{min}} = 24$  of the last 30 CPU usage values exceed 85%.

**Log rule (multiple mining-related clues).** Let the set of *distinct* matched keywords be

$$\mathcal{S}(L) = \{k \in \mathcal{K} \mid \exists \ell \in L \text{ such that } k \preceq_{\text{ci}} \ell\}, \quad C_{\text{log}}(L) = |\mathcal{S}(L)|. \quad (4)$$

Then the log-side condition passes iff at least  $s_{\text{min}} = 6$  distinct mining clues appear:

$$\text{VerifyLog}_{\text{mining}}(L) = \mathbf{1}[C_{\text{log}}(L) \geq s_{\text{min}}]. \quad (5)$$

**Final symbolic verification for mining.** The mining scenario is accepted iff both metric and log checks pass:

$$\text{Verifier}_{\text{mining}}(M, L) = \begin{cases} 1, & \text{if } \text{VerifyMetric}_{\text{mining}}(M) = 1 \wedge \text{VerifyLog}_{\text{mining}}(L) = 1, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$