

LabelEase: A Semi-Automatic Tool for Efficient and Accurate Trace Labeling in Microservices

Shenglin Zhang¹, Zeyu Che¹, Zhongjie Pan¹, Xiaohui Nie², Yongqian Sun¹, Lemeng Pan³, Dan Pei⁴

Nankai University¹, CNIC², Huawei³, Tsinghua University⁴















Trace and Span





An Example of Span Log

- Large-scale microservice systems could contain tens of thousands of service instances
- Trace data is crucial for understanding system behavior, diagnosing performance issues, and ensuring reliability

Trace-Based Anomaly Detection and Root Cause Localization





- Anomaly detection identifies anomalies in service execution and generates alerts
- Root cause localization identifies the specific component or subsystem
 responsible for each failure

Why Traces Need to Be Labeled



- Trace-based anomaly detection methods usually apply deep learning methods to learn the normal patterns from "normal data", the high-quality data without anomalies are crucial
- High-quality data is important to evaluate and improve the performance of trace-based anomaly detection and root cause localization methods
- There are only a few public trace datasets
 - Limited number of traces
 - Limited types of failures
 - Not generalized



Manual Trace Labeling Is Time-Consuming and Labor-Intensive







Our goal is to design a semi-automatic trace labeling tool for efficient and accurate trace labeling

Content





Implementation

Conclusion

Massive Trace Data





- Microservice systems could generate a massive amount of trace data every day for labeling
- E.g., the microservice systems of eBay generate about 150 billion traces everyday!

Complex Trace Structure



- A trace is a complex structure arising from the hierarchy of parallel service invocations
- It can contain dozens to hundreds of services



Hard-To-Determine Root Cause





- Root cause labeling needs to consider multiple traces in a period
- It is challenging to check all the anomalous trace data and conclude the root cause information



LabelEase

A semi-automatic trace labeling tool for efficient and accurate trace labeling

11







Overview





Graph-Based Trace Representation



- Data Preprocessing
- GNN Training
- Vectorized Representation



Graph-Based Trace Representation



- Data Preprocessing
 - Convert the raw trace data into a graph
 - Focus on semantic information, temporal features, status codes within the span, and the calling relationships between different services
- GNN Training
 - Graph contrastive learning
 - Differentiate between various traces based on their representations
- Representation vectorization



Anomaly Labeling



• Hybrid Representative Selection

- Partition traces into different classes
- Select the cluster center point as the most representative trace for operators to label
- Human Feedback
 - Operators label the most representative trace in each cluster
 - LabelEase automatically labels the remaining traces



Anomaly Labeling



 Use active learning to select the most informative data points from an unlabeled dataset for labeling through human feedback

 Minimize the labeling effort while ensuring high-quality data labels









Trace Aggregation

Partition the trace data into abnormal periods

- ≻Root Cause Localization
 - > The spectrum-based fault localization (SBFL) technique
 - Service dependency topology
 - Service latency information











After LabelEase selects representative traces for labeling, operators will label the given traces on this page



Trace Anomalies Labeling



Trace and

span-

related

information

• Span ID, start time, serice name, operation name, status code, and duration



Trace Anomalies Labeling



• The average value, standard deviation, and other information about each similar call path in the entire dataset as a labeling



Trace Anomalies Labeling



• Each trace is visually distinguished: green signifies a label of "normal", red indicates a label of "abnormal", and black denotes traces awaiting







After all traces are labeled, operators will go to the page of labeling root cause.





- LabelEase identifies the impact of each fault by merging the anomalous traces occurred closely
- Red means it has been labeled and black denotes period awaiting labeling





- For each period, a microservice system relationship diagram is featured in the center of the page
- The higher the suspicious score, the more likely the root cause is, and the darker the service will be displayed on the interface





- Upon selecting a service, detailed information about it is displayed
- Operators can set the type of fault to label the root cause
- The service labeled as the root cause is highlighted in red







Datasets and Metrics



• We use two datasets, denoted as Dataset 1 and Dataset 2.

	Source	#Normal traces	#Anomalous traces
Dataset 1	A benchmark microservice system developed by us	103078	27285
Dataset 2	A top-tier commercial bank	93579	19207

• Evaluation metrics

$$Precision = \frac{TP}{TP + FP} \qquad Recall = \frac{TP}{TP + FN}$$

$$F_1 score = 2 \times \frac{Precision \times Rcall}{Precision + Recall}$$



- LabelEase shows superior performance over all baseline approaches across both datasets, achieving F1-scores of 0.99 and 0.98, respectively
- It shows superior efficiency compared to the baseline methods, boasting the shortest time overhead and highest labeling efficiency

Approach	$\mathscr{D}1$				$\mathscr{D}2$			
Approach	Precision	Recall	F_1 -score	Time	Precision	Recall	F_1 -score	Time
LabelEase	1	0.98	0.99	6.53s	0.96	0.99	0.98	21.68s
MultimodalTrace [37]	0.2	0.15	0.17	1.9min	0.17	0.15	0.16	1.6min
TraceAnomaly [8]	0.94	0.67	0.78	27.2min	0.21	0.2	0.2	22.2min
CRISP [9]	0.8	0.57	0.67	26.1min	0.24	0.21	0.23	17.3min
TraceCRL [10]	0.39	0.28	0.33	8.9h	0.48	0.43	0.45	7.2h
TraceVAE [11]	0.4	0.3	0.35	1.7h	0.14	0.77	0.23	51.3min
TraceSieve [3]	1	0.74	0.85	7.6min	0.17	0.15	0.15	9.8min

THE EFFECTS OF LabelEase IN COMPARISON WITH DIFFERENT APPROACHES ON TWO DATASETS

TABLE I

Effect of Graph-based Trace Representation



 LabelEase with the graph-based representation outperforms other trace vectorization methods





• The K-means algorithm shows superior effects and the shortest computational time compared with baseline methods

TABLE II THE EFFECTS OF THE DIFFERENT CLUSTER METHODS ON TWO DATASETS

Approach	$\mathscr{D}1$			$\mathscr{D}2$				
Approach	Precision	Recall	F_1 -score	Time(s)	Precision	Recall	F_1 -score	Time(s)
using K-means	1	0.98	0.99	6.53	0.96	0.99	0.98	21.68
using hierarchical clustering	0.99	0.98	0.99	96.9	0.96	0.95	0.96	51.98
using DBSCAN	0.99	0.92	0.95	116.07	0.83	1	0.91	97.55
using Spectral clustering	0.99	0.98	0.98	200.84	0.96	0.93	0.94	264.24
using random selection	0.8	0.2	0.32	-	0.95	0.33	0.49	-

Sensitivity of the Number of Traces to be Labeled



 How the F1-score of LabelEase changes with different trace numbers of labels





• The precision of fault period detection is 0.89.

- As long as an anomalous trace is labeled, it is possible to accurately identify the period in which the corresponding fault occurred, thus enabling precise root cause labeling
- The recall is 0.44 when localizing the fault period
 - Anomalies in the unreported fault periods are not reflected in the traces but primarily in metrics and logs











- A novel tool LabelEase, a semi-automatic trace labeling tool for efficient and accurate trace labeling
- A series of experimental studies to evaluate LabelEase's effectiveness and efficiency using two datasets
- Publish a high-quality dataset
 <u>https://doi.org/10.5281/zenodo.13338156</u>
- LabelEase will surely contribute to the rapid development of AIOps



Thank you!