

MicroRCA-Agent: Microservice Root Cause Analysis Method Based on Large Language Model Agents

Pan Tang¹, Shixiang Tang¹, Huanqi Pu¹, Zhiqing Miao², Zhixing Wang³

¹School of Communication and Information Engineering, Shanghai University, Shanghai, China

²School of Communication and Electronic Engineering, East China Normal University, Shanghai, China

³School of Information and Electronics, Beijing Institute of Technology, Beijing, China

ptang@shu.edu.cn, tangsx624@163.com, 1287056469@qq.com, 52275904009@stu.ecnu.edu.cn, wangwangzx@bit.edu.cn

Abstract

This paper presents MicroRCA-Agent, an innovative solution for microservice root cause analysis based on large language model agents, which constructs an intelligent fault root cause localization system with multimodal data fusion¹. The technical innovations are embodied in three key aspects: First, we combine the pre-trained Drain log parsing algorithm with multi-level data filtering mechanism to efficiently compress massive logs into high-quality fault features. Second, we employ a dual anomaly detection approach that integrates Isolation Forest unsupervised learning algorithms with status code validation to achieve comprehensive trace anomaly identification. Third, we design a statistical symmetry ratio filtering mechanism coupled with a two-stage LLM analysis strategy to enable full-stack phenomenon summarization across node-service-pod hierarchies. The multimodal root cause analysis module leverages carefully designed cross-modal prompts to deeply integrate multimodal anomaly information, fully exploiting the cross-modal understanding and logical reasoning capabilities of large language models to generate structured analysis results encompassing fault components, root cause descriptions, and reasoning trace. Comprehensive ablation studies validate the complementary value of each modal data and the effectiveness of the system architecture. The proposed solution demonstrates superior performance in complex microservice fault scenarios, achieving a final score of 50.71. The code has been released at: <https://github.com/tangpan360/MicroRCA-Agent>.

Overall Design

The proposed solution adopts a modular architectural design, with the overall system comprising five core modules: data preprocessing module, log fault extraction module, trace anomaly detection module, metric fault summarization module, and multimodal root cause analysis module. The modules are designed with loose coupling, facilitating data interaction through function encapsulation, which ensures both system integrity and module independence with enhanced scalability.

The data preprocessing module undertakes the critical responsibility of system data standardization. Its core functionalities include: parsing fault time period information from input.json,

¹This work is a technical report of our solution at the 2025 (8th) CCF International AIOps Challenge. The official website of the Challenge is <https://challenge.aiops.cn>.

and standardizing timestamp formats across log, trace, and metric modalities to ensure accuracy in subsequent cross-modal data correlation.

The log fault extraction module is based on a "template + rule" processing architecture. It first designs multiple regular expression rules covering different fault types. The module utilizes the pre-trained Drain model(He et al., 2017) (error.log-drain.pkl) for automatic log template extraction, merging log entries with the same semantics into templates with their occurrence counts through template recognition and deduplication mechanisms, with precise control over template matching parameters via the drain3.ini configuration file. The module implements a multi-level data filtering mechanism: file localization → time window filtering → error keyword filtering → extraction and reconstruction of core fields → fault template matching and standardization → sample deduplication and frequency statistics → name mapping and service extraction, effectively compressing massive raw logs into high-quality fault features.

The trace fault detection module adopts a hybrid detection strategy of "unsupervised learning + rule enhancement". At its core, it leverages the Isolation Forest(Liu et al., 2008) algorithm to construct anomaly detection model, which uses duration data from parent_pod→child_pod call chains during normal periods for model training. The model training process involves extracting normal data from 50 random samples, using 30-second sliding windows to extract temporal features, and generating the trace_detectors.pkl model file upon completion. During the detection phase, anomaly detection is performed on the duration of parent_pod→child_pod call chains in fault periods to filter out significantly anomalous call chains and their corresponding durations. Additionally, it implements a status analysis function, extracting status.code and status.message information from traces to identify call failure patterns. Finally, it outputs two sets of results respectively: statistical analysis of the top 20 most frequent anomalous call combinations, and detailed status anomaly information for the top 20 most frequent cases.

The metric fault summarization module designs a statistical symmetric ratio filtering strategy combined with a hierarchical two-stage LLM analysis approach. First, it performs mathematical statistical analysis on metric indicators and uses symmetric ratios to determine whether the statistical indicators of normal intervals and fault intervals are highly similar, thereby filtering normal data and significantly reducing the large model context length, enabling it to focus more on potentially anomalous data. After filtering, it conducts two-stage LLM analysis: the first stage is based on Application Performance Monitoring (APM) data and database component (TiDB) data, performing phenomenon summarization at both service level and pod level, including business metrics such

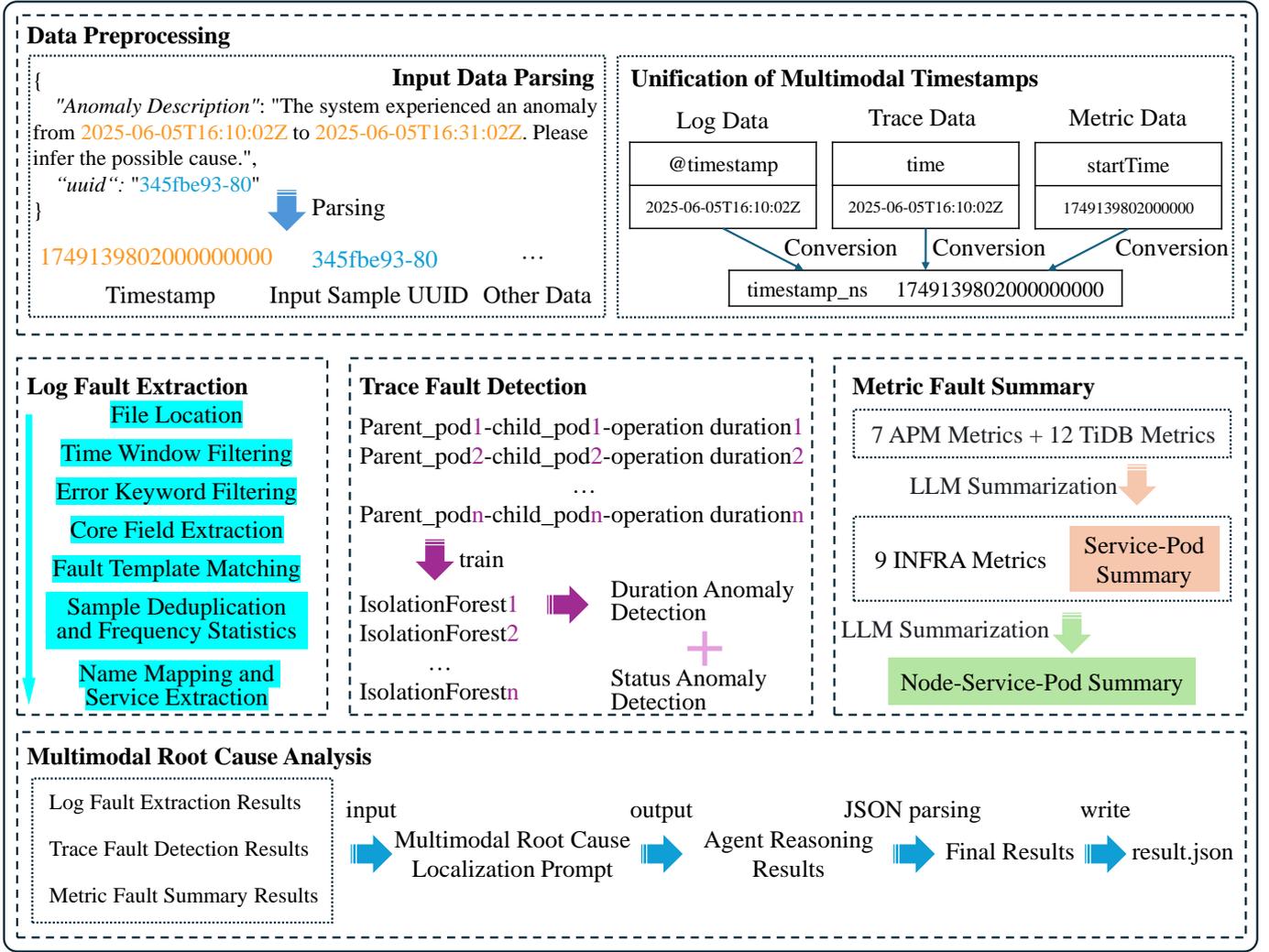


Figure 1: Overall architecture of MicroRCA-Agent.

as microservice request response and anomaly ratios, as well as related indicators of database auxiliary components TiDB, TiKV, and PD. The second stage combines the business metric phenomenon summary from the first stage with node-level infrastructure data to conduct comprehensive phenomenon analysis across three levels: service, pod, and node. This approach effectively controls computational costs while ensuring analytical depth.

The multimodal root cause analysis module is designed with specialized multimodal prompts at its core, supporting flexible combinations of log, trace, and metric data across three modalities. The module adopts a multi-processing strategy to improve system processing efficiency and integrates comprehensive fault tolerance mechanisms and retry strategies. It ultimately outputs structured root cause analysis results containing component, reason, and reasoning trace, achieving a complete closed loop from phenomenon observation to root cause inference.

2 Innovation and Practicality

The innovation of this solution is embodied in its multi-level feature extraction and reasoning architecture design. First, the log, trace, and metric modules each adopt highly targeted processing and filtering strategies: the Drain algorithm is combined with

multi-level data filtering to compress log redundancy; Isolation Forest and status code checking complement each other to identify trace anomalies; statistical significance-based filtering is applied to extract high-value metric features. This modality-specific "precise extraction-information compression" mechanism, tailored to each data type (log, trace, metric), significantly reduces the computational and contextual overhead of subsequent analysis while ensuring the integrity of key information. Second, the metric analysis component incorporates a statistical symmetric ratio filtering mechanism and a two-stage phenomenon summarization framework. After filtering out potentially anomalous data, the framework organically integrates service/instance and container/node perspectives, while preserving the correlation between business operations and infrastructure at the phenomenon level, thereby providing semantically interpretable evidence directly usable for subsequent root cause localization.

In terms of practicality, this solution emphasizes a modular and scalable architectural design, enabling flexible adaptation to different business and system environments. The functional modules (log extraction, trace detection, metric summarization, root cause analysis) are decoupled from one another, allowing for independent deployment or on-demand combination. This design facilitates rapid implementation in microservice architectures of varying

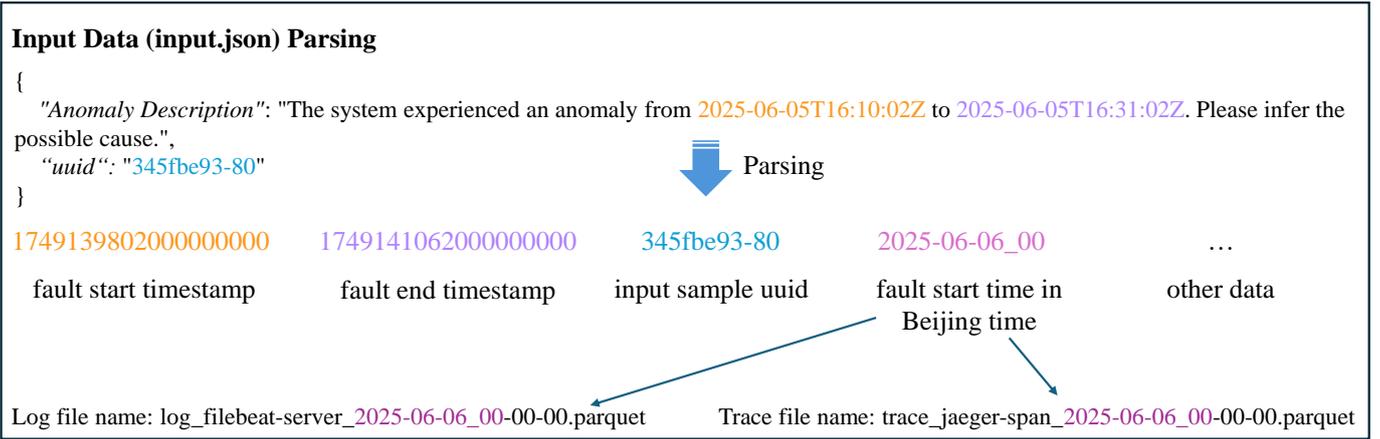


Figure 2: Input data parsing.

```

uuid, start_time_beijing, end_time_beijing, start_timestamp, end_timestamp, start_time_hour
345fbe93-80,2025-06-06_00-10-02,2025-06-06_00-31-02,1749139802000000000,1749141062000000000,2025-06-06_00

```

Figure 3: Preprocessed input data.

scales and complexities. In terms of processing efficiency, the solution significantly reduces the amount of irrelevant data entering the analysis phase via strategies like hierarchical filtering. This reduction in irrelevant data lowers computational and storage overhead, thus making the solution well-suited for high-concurrency, large-data-volume production environments.

3 Detailed Design

Centering on the demand for fault root cause localization in microservice systems, this solution constructs a complete system architecture consisting of five modules: data preprocessing, log fault extraction, trace fault detection, metric fault summarization, and multimodal root cause analysis (as shown in Figure 1). First, the data preprocessing module realizes the structured parsing of input fault information and the unification of multimodal timestamps. Subsequently, the log module performs multi-level data filtering on large-scale logs and uses the Drain algorithm to extract log templates for deduplication; the trace module adopts a dual strategy combining duration performance analysis and status verification to identify trace anomalies; the metric module leverages large language models to conduct phenomenon summarization for multi-level monitoring metrics. Finally, all modal data are aggregated into the multimodal root cause analysis module, which utilizes the cross-modal reasoning capability of large language models to achieve the localization of faulty components, the explanation of causes, and the output of evidence chains.

3.1 Data Preprocessing

The data preprocessing module focuses on two core functions: parsing input data and uniformly processing multimodal timestamps. It ensures the consistency and accuracy of subsequent processing for each modal data, thereby laying a data foundation for the entire fault root cause localization system.

3.1.1 Design of Input Data Parsing

For input data, a regex-based extraction strategy is adopted to perform structured processing on raw fault description files. The system accepts input data in JSON format, which contains two core components: an anomaly description and a unique identifier (uuid). Among these components, the most critical pieces of information are the fault start/end times (included in the anomaly description) and the uuid.

To facilitate accurate matching of subsequent multimodal data and filtering of data within the fault time window, the system first parses and extracts key information from the input files. The core information extracted includes the unique fault identifier and the fault start/end times, as illustrated in Figure 2.

Timestamp Extraction Mechanism: The system adopts the ISO 8601 time format standard and implements automated timestamp identification through the regular expression pattern $(\backslash\{4\}-\backslash\{2\}-\backslash\{2\}T\backslash\{2\}:\backslash\{2\}:\backslash\{2\}Z)$. The extraction rule follows the strategy where the first match is taken as the fault start time and the second match as the fault end time.

Time Index Generation: To enhance the efficiency of subsequent data correlation, the system constructs a time indexing mechanism. It generates time identifiers in the format "year-month-day_hour" (e.g., 2025-06-06.00) for rapid localization of corresponding data files. Additionally, fault times are converted into 19-digit nanosecond-level timestamps, providing a high-precision reference for accurate time-range filtering. The preprocessed input data is specifically illustrated in Figure 3.

3.1.2 Unification of Multimodal Data Timestamps

In response to the distinct format characteristics of the three data types (log, trace, and metric), the system is designed with differentiated timestamp unification strategies in a systematic manner. These strategies enable the standardization of cross-modal time references, as specifically illustrated in Figure 4.

Timestamp Standardization for Log Data: For log data, the @timestamp field in ISO 8601 format is adopted as the time ref-

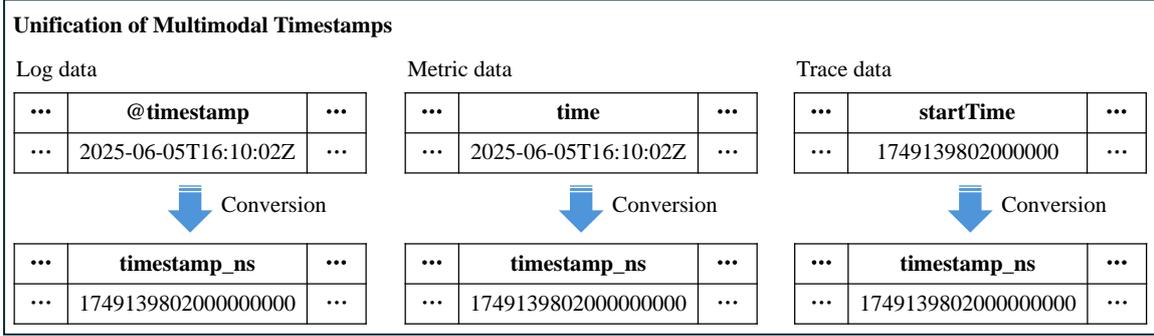


Figure 4: Unification of multimodal data timestamps.

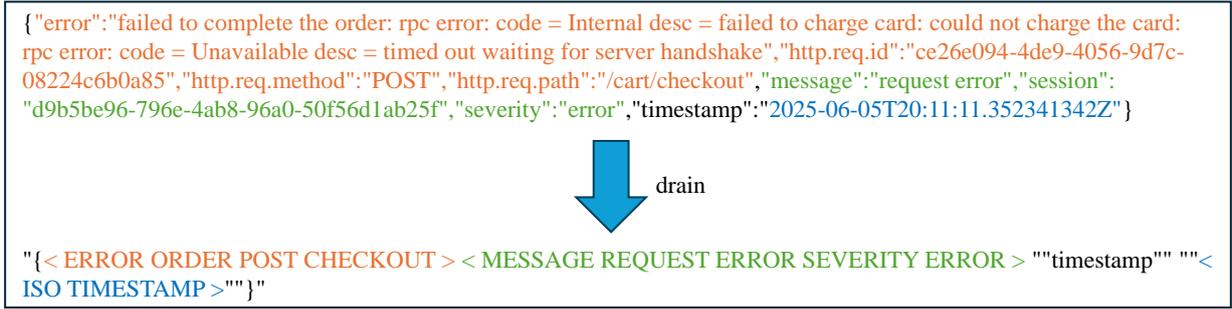


Figure 5: Log template extraction using Drain.

erence. The system parses the time format to convert raw timestamps into unified 19-digit nanosecond-level timestamps. After processing, the log data is sorted in ascending order of timestamps, ensuring the temporal consistency of the data and the accuracy of subsequent analyses.

Timestamp Standardization for Metric Data: Metric data stores time information in the "time" field, which also adheres to the ISO 8601 format standard. Given the multi-level distributed nature of metric data in the storage architecture, the system employs a recursive search strategy to traverse all subdirectory structures, ensuring complete coverage of metric files stored in a decentralized manner. The timestamp conversion logic remains consistent with that applied to log data.

Timestamp Standardization for Trace Data: Trace data features a unique time storage format, with microsecond-level timestamps stored in the "startTime" field. In this study, precision conversion is applied to extend microsecond-level timestamps to nanosecond-level ones (by multiplying by a factor of 1000), achieving time precision alignment with other modal data.

3.2 Log Fault Extraction

The log fault extraction module serves as a core component of the fault root cause localization system, employing a multi-level filtering strategy to efficiently convert massive volumes of raw logs into structured fault feature information. The fault information compression functionality of this module is based on a pre-trained Drain algorithm model. It achieves effective compression of fault logs by extracting templates from logs containing error fields and performing deduplication on logs sharing the same template, thereby providing high-quality structured log data input for subsequent multimodal fault analysis.

3.2.1 Log Parsing Principles and Model Construction

a Challenges in Unstructured Log Processing

As unstructured information, logs record system operational behaviors in detail. However, they also contain a large number of irrelevant variables (e.g., timestamps, IP addresses, port numbers) that impair analysis efficiency. While these variables hold significance in logs, they tend to result in a multitude of duplicate records with identical semantics but different expressions during fault pattern recognition. Such duplicate records occupy valuable context space in large language models without providing additional valid information.

b Selection and Application of Drain Algorithm

After a comprehensive evaluation of existing log parsing techniques, this study adopts the widely recognized Drain algorithm as the core log parser. By constructing a parsing tree, the Drain algorithm can effectively identify fixed and variable parts within logs, and categorize logs with identical semantic patterns into the same template. This process enables the conversion of unstructured logs into structured features while compressing a large amount of redundant logs. Specifically, the Drain algorithm extracts log templates, removes irrelevant variables (e.g., timestamps, IDs), and merges logs sharing the same template into a single entry, thereby reducing the data volume effectively. Figure 5 illustrates the log template extraction process: in its output, all variables are parsed into invariant content by Drain, which mitigates the interference caused by irrelevant variables.

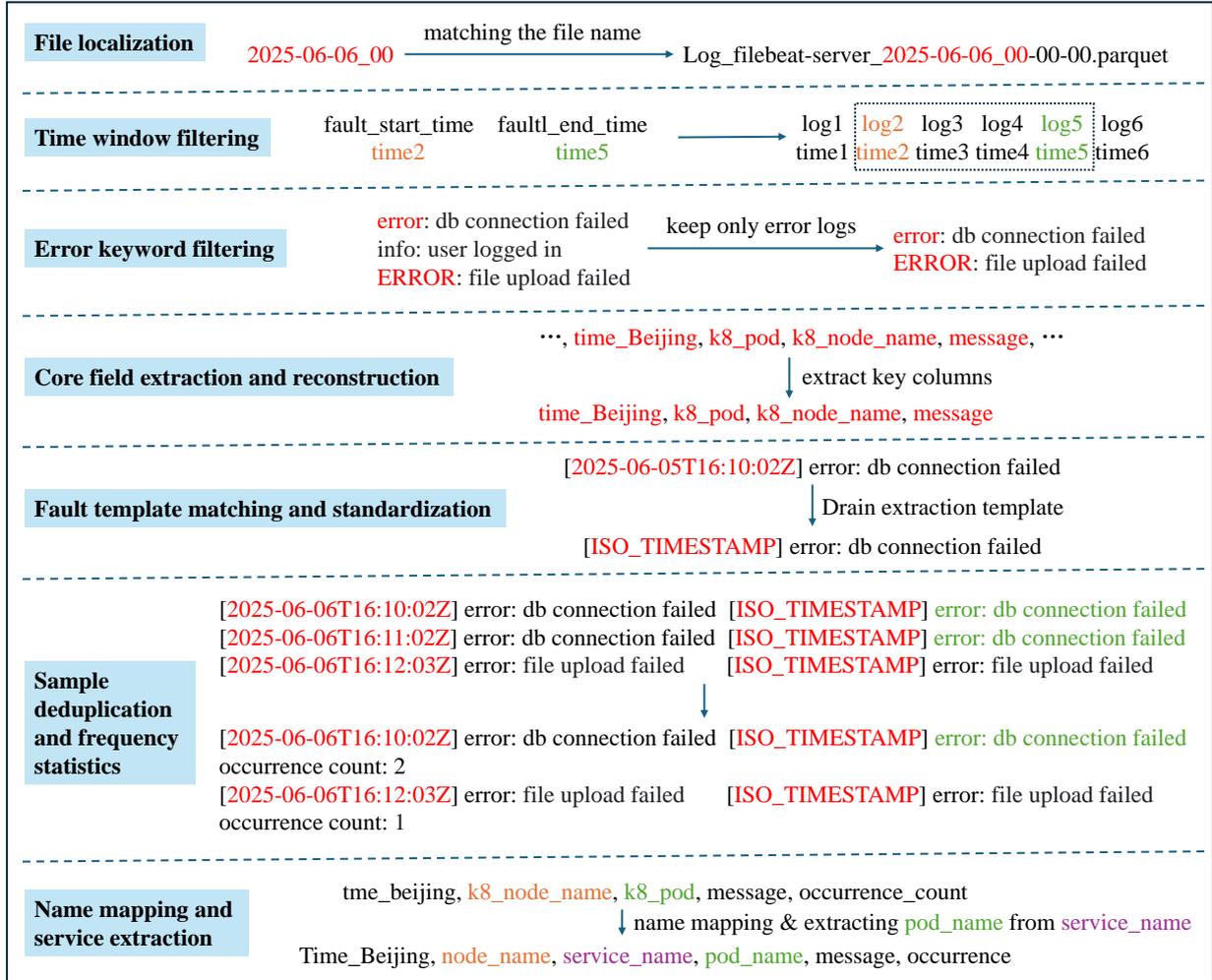


Figure 6: Multi-level data filtering and processing workflow.

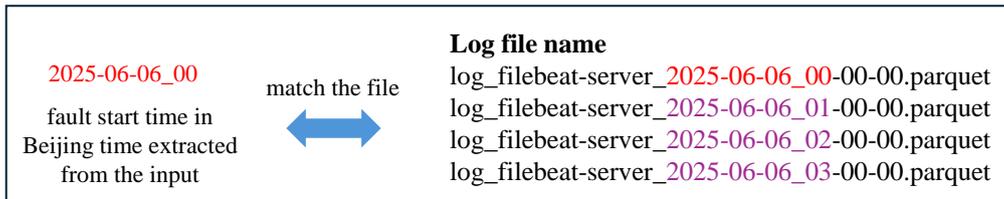


Figure 7: Matching between time in input and log file names.

c Construction of the Pre-trained Model

The system filters all log records containing the "error" field from the full-volume log data in phaseone² to use as training samples, thereby constructing a pre-trained Drain model. Through the learning and analysis of massive error logs in the microservice environment, the model successfully extracts and builds 156 representative fault templates, which cover the main fault pattern types in the microservice system in a relatively comprehensive manner.

3.2.2 Multi-level Data Filtering and Processing Workflow

The multi-level data filtering and processing workflow is specifically illustrated in Figure 6, proceeding from file localization,

²The dataset consists of two parts, phaseone and phasetwo; here we use the phaseone data for training. Source: <https://www.aiops.cn/gitlab/aiops-live-benchmark/phaseone/>.

through time window filtering, error keyword filtering, extraction and reconstruction of core fields, fault template matching and standardization, sample deduplication and frequency statistics, to name mapping and service extraction. This process ultimately converts massive volumes of logs into structured fault features while achieving effective content compression.

a File Localization

The system performs file matching using time identifiers in the "year-month-day_hour" format, based on the time information of each input item output from the preprocessing stage, as illustrated in Figure 7. By searching for log files corresponding to the target time period within the project data directory, the system ensures accurate localization of data sources within the fault time window.

node_name	service_name	pod_name	message	occurrence_count
aiops-k8s-04	frontend	frontend-2	“{“error”：“could not retrieve product: ...”	occurrence_count: 81
aiops-k8s-03	frontend	frontend-0	“{“error”：“could not retrieve product: ...”	occurrence_count: 9
...				
aiops-k8s-04	frontend	frontend-2	“{“error”：“failed to get product ...”	occurrence_count: 1

Figure 8: Output of log fault extraction.

b Time Window Filtering

Based on the nanosecond-precision timestamps of the fault start and end times, strict time-boundary filtering is performed on log data. The system conducts a range query using the timestamp field, ensuring that only log records within the fault time period are retained and effectively eliminating interference from irrelevant information in the time.

c Error Keyword Filtering

A filtering rule based on the "error" keyword is implemented for log message fields. After analyzing a large volume of log data, it is observed that faults are likely to occur in a large number of log entries accompanied by the "error" field. Therefore, log entries containing "error" information are automatically identified and extracted, while logs related to normal business operations are filtered out. This significantly enhances the relevance and efficiency of subsequent analyses.

d Extraction and Reconstruction of Core Fields

Key information essential for analysis is accurately extracted from multiple fields of raw logs, including critical dimensions such as time information (time_beijing), container identifier (k8_pod), node information (k8_node_name), and error message content (message). A field simplification strategy is employed to reduce data redundancy while ensuring information integrity.

e Fault Template Matching and Standardization

The pre-trained Drain model is utilized to perform template matching for each error log. The system converts raw log messages into standardized template representations, removes irrelevant variables, and retains only core fault semantic content, thereby enabling the extraction from individual logs to pattern categories.

f Sample Deduplication and Frequency Statistics

A deduplication strategy is applied to duplicate log records with the same container and template combination (e.g., [[2025-06-06 08:00, frontend-0, template0, message0], [2025-06-06 09:00, frontend-0, template0, message1]]). The system retains the first occurrence record of each combination and simultaneously counts the occurrence frequency of the combination (e.g., [2025-06-06 08:00, frontend-0, message0, Occurrence Count: 2]). This quantitative labeling provides an objective basis for evaluating the severity of faults.

g Name Mapping and Service Extraction

The k8_node_name and k8_pod fields in logs are mapped and converted into node_name and pod_name respectively. Corresponding service information is then extracted by parsing the pod_name (e.g., frontend-0 → frontend). Finally, the system reconstructs the data into a standardized multi-level format, which includes dimensions such as node, service, pod, error message, and frequency statistics. This provides a structured data foundation for subsequent multimodal fault analysis.

Through the processing workflow described above, the log fault extraction module can efficiently compress massive volumes of raw log data into a high-quality set of structured fault features. Specific examples of the finally extracted faults are illustrated in Figure 8. This module not only significantly reduces the complexity of data processing but also maintains the integrity and accuracy of fault information, thereby providing log-level data support for subsequent agent-based multimodal fault root cause localization.

3.3 Trace Fault Detection

The trace fault detection module is a critical component for fault root cause localization in microservice systems. It adopts a dual anomaly detection strategy to identify abnormal patterns in microservice traces from both performance and status dimensions. Based on the Isolation Forest algorithm and a direct status code inspection method, this module effectively detects duration anomalies and status anomalies. The two types of detection results are simultaneously fed into multimodal analysis, providing high-quality structured trace data input for subsequent multimodal fault analysis.

3.3.1 Principle of Dual Anomaly Detection and Model Construction

Distributed traces in microservice architectures contain abundant performance and status information, yet abnormal patterns exhibit diverse characteristics. Performance anomalies are primarily manifested as a significant deviation of trace duration from the normal range, which requires the establishment of a normal behavior baseline through machine learning methods. In contrast, status anomalies are characterized by the occurrence of explicit error codes and can be identified via direct inspection. The combination of these two anomaly detection mechanisms enables the provision of more comprehensive fault feature information.

a Principle of Isolation Forest for Performance Anomaly Detection

For detecting duration-based performance anomalies, the Isolation Forest unsupervised learning algorithm is employed. This algorithm operates on the core assumption that "anomalies are more

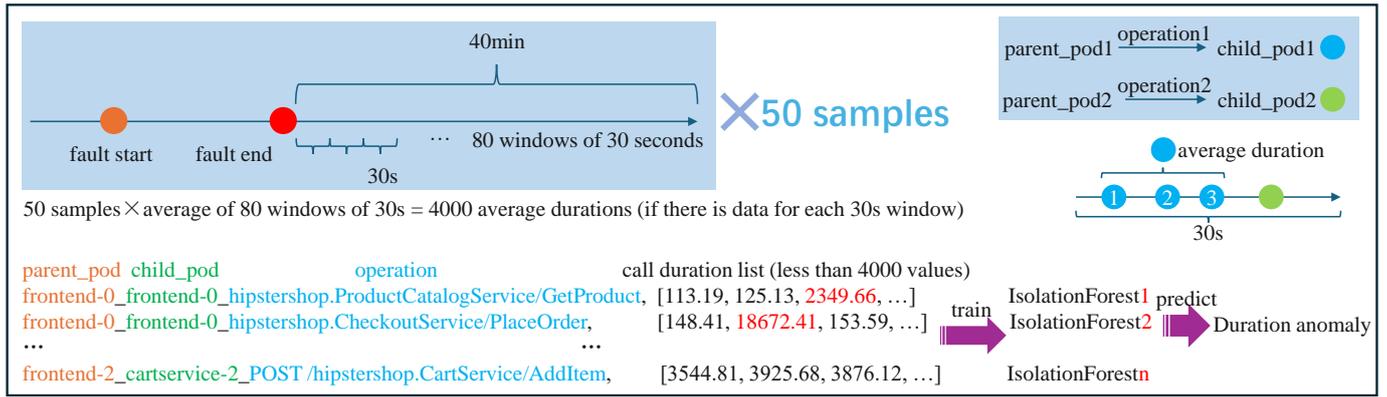


Figure 9: Construction of Isolation Forest training data and model training.

easily isolated in the feature space” and quantifies the anomaly degree of data points by constructing multiple randomly partitioned trees. Isolation Forest is particularly suitable for anomaly detection in continuous numerical data, enabling automatic identification of significant deviation patterns in duration data without the need for labeled anomaly samples.

b Direct Inspection Mechanism for Status Codes

A direct inspection mechanism based on status codes is constructed, eliminating the need for a training process. The system directly analyzes the status.code and status.message fields within the tags of trace data, and filters out all abnormally statused invocations through a simple conditional judgment ($\text{status.code} \neq 0$). This method can directly capture explicit abnormal states and specific information such as system errors, timeouts, and connection failures, providing deterministic anomaly evidence for fault root cause localization.

c Construction of Isolation Forest Training Data and Model Training

As illustrated in Figure 9, the construction of Isolation Forest training data and the subsequent model training primarily consist of the following components:

Training Data Source: In this scheme, data from the phase one stage is used to construct the training set. Fifty fault samples are randomly selected, and a 40-minute time window immediately after the end of each fault is extracted as “normal period” data. This approach is based on the assumption that system behavior tends to return to normal after fault recovery, ensuring that the majority of the training data consists of normal data.

Grouped Training Strategy: Training data is grouped according to the service invocation combination of “parent_pod - child_pod - operation name”, and an independent Isolation Forest detector is trained for each unique service invocation pattern. This fine-grained grouping strategy accounts for differences in performance baselines across various service invocations, thereby improving the accuracy of anomaly detection.

Sliding Window Averaging Processing: During the training phase, a 30-second sliding window is employed to process the duration data as an event sequence. Since multiple invocation combinations of “parent_pod - child_pod - operation name” may occur

within a single 30-second sliding window, and each combination may appear one or more times, the average duration is calculated for each combination that occurs at least once within the window. If a combination does not appear, it is recorded as “None”. This processing method effectively eliminates random fluctuations in individual invocation samples, stabilizes the duration feature, and thereby more accurately reflects the actual performance level of service invocations.

Anomaly Detection Threshold Setting: For the Isolation Forest model, the contamination parameter is set to 0.01 (indicating that the expected proportion of abnormal samples to the total is 1%), and $n_estimators = 100$ is used to construct 100 decision trees. The model outputs an anomaly score: a score of -1 indicates an abnormal sample, while a score of 1 indicates a normal sample.

3.3.2 Multi-level Trace Data Processing Workflow

a File Localization and Time Filtering

The system performs precise matching of trace files in the “year-month-day_hour” format based on the time information of each input parsed in the preprocessing stage. Using nanosecond-level timestamps of the fault’s start and end times, it executes time window filtering on the trace data, ensuring that the analysis scope covers the complete fault time period.

b Invocation Relationship Mapping and Structured Extraction

Key dimensional information is extracted from complex trace structures, including pod_name, service_name, and node_name obtained by parsing the process field, as well as the establishment of complete parent-child relationship mappings for call chains using spanID and references. The system converts unstructured trace data into invocation representations containing “parent_pod - child_pod - operation name”.

c Duration Anomaly Detection Processing Flow

Prediction Data Processing: Anomaly detection prediction is performed on all trace data within the fault time period, employing the same 30-second sliding window strategy as in the training phase to calculate the average duration within each time window. This consistent processing ensures the consistency of feature distribution between the training and prediction phases.

node_name,	service_name,	parent_pod,	child_pod,	operation_name,	normal_avg_duration,	anomaly_avg_duration,	anomaly_count
aiops-k8s-07,	frontend,	frontend-1,	frontend-1,	histershop.ProductCatalogService/GetProduct,	11203.63505,	56375295.28801,	occurrence_count: 32
aiops-k8s-04,	frontend,	frontend-2,	frontend-2,	histershop.ProductCatalogService/GetProduct,	11480.30423,	57508444.49971,	occurrence_count: 31
...							
aiops-k8s-03,	frontend,	frontend-0,	frontend-0,	histershop.ProductCatalogService/GetProduct,	11371.84608,	56044805.23074,	occurrence_count: 5

Figure 10: Example output of anomalous traces detected by Isolation Forest.

node_name,	service_name,	parent_pod,	child_pod,	operation_name,	status_code,	status_message,	occurrence_count
aiops-k8s-04,	frontend,	N/A,	frontend-2,	histershop.Frontend/Recv.,	13,	HTTP status cde: 500,	occurrence_count: 161
aiops-k8s-03,	frontend,	N/A,	frontend-0,	histershop.Frontend/Recv.,	13,	HTTP status cde: 500,	occurrence_count: 159
...							
aiops-k8s-07,	frontend,	frontend-1,	frontend-1,	histershop. ProductCatalogService/GetProduct,	1,	context canceled,	occurrence_count: 131

Figure 11: Example output of anomalous traces detected by status code.

Isolation Forest Anomaly Prediction: The pre-trained detection model is used to identify anomalies in the duration features during the fault period. A dedicated detector corresponding to each service invocation combination is employed for prediction (if an invocation combination that was not trained before exists during the fault period, it is ignored and no prediction is performed), and anomaly labels are output (with -1 indicating an anomaly and 1 indicating normal).

Duration Anomaly Result Output: The system outputs the top 20 most frequent duration anomaly combinations, which include the following specific details:

- node_nameNode: identifier where the anomaly occurs;
- service_name: Name of the anomalous service;
- parent_pod: Identifier of the caller container;
- child_pod: Identifier of the callee container;
- operation_name: Name of the specific operation;
- normal_avg_duration: Average invocation latency during the normal period (derived from training data statistics);
- anomaly_avg_duration: Average invocation latency during the anomalous period (derived from actual statistics of the fault phase);
- anomaly_count: Frequency of anomaly occurrences (format: Number of occurrences: N);

A specific example of the output of anomalous traces detected by the Isolation Forest is shown in Figure 10.

d Status Anomaly Detection Processing Flow

Direct Status Check: Status code checks are performed on all trace data within the fault time period, with no training process required. Information on status.code and status.message is parsed from the tags field of traces, and all anomalous status invocations are directly identified through conditional filtering (status.code \neq 0).

Anomalous Status Pattern Statistics: The identified anomalous statuses are grouped and counted by service invocation combinations, and the occurrence frequency of each anomalous pattern is calculated.

Status Anomaly Result Output: The system outputs the top 20 most frequent status anomaly combinations, which include the following specific details:

- Node_name: Node identifier where the anomaly occurs;
- Service_name: Name of the anomalous service (Redis is automatically mapped to redis-cart);
- Parent_pod: Identifier of the caller container;
- Child_pod: Identifier of the callee container;
- Operation_name: Name of the specific operation;
- Status_code: Specific error status code;
- Status_message: Detailed description of the error status;
- Occurrence_count: Frequency of status anomaly occurrences (format: Number of occurrences: N);

A specific example of the output of anomalous traces detected by the Status Check is shown in Figure 11.

Through the aforementioned dual anomaly detection workflow, the trace fault detection module achieves an organic combination of machine learning-based anomaly identification and rule-based direct anomaly checking. Duration anomaly detection provides quantitative performance anomaly analysis via the Isolation Forest algorithm, while status anomaly detection enables definitive error status identification through direct checks. The top 20 most critical anomaly patterns output by each of the two detection strategies collectively provide a trace data foundation for subsequent fault root cause localization.

3.4 Metric Fault Summary

The Metric Fault Summary Module serves as a full-stack monitoring and analysis component for faults in distributed microservice systems. It adopts a two-stage phenomenon summary strategy based on large language models (see Figure 12). By leveraging rule-based filtering and statistical comparison methods, it identifies significantly anomalous metrics. Furthermore, it utilizes the reasoning capabilities of large language models to perform phenomenon induction and pattern recognition on complex multi-dimensional monitoring data, thereby providing high-quality phenomenon description inputs for subsequent multimodal fault root cause analysis.

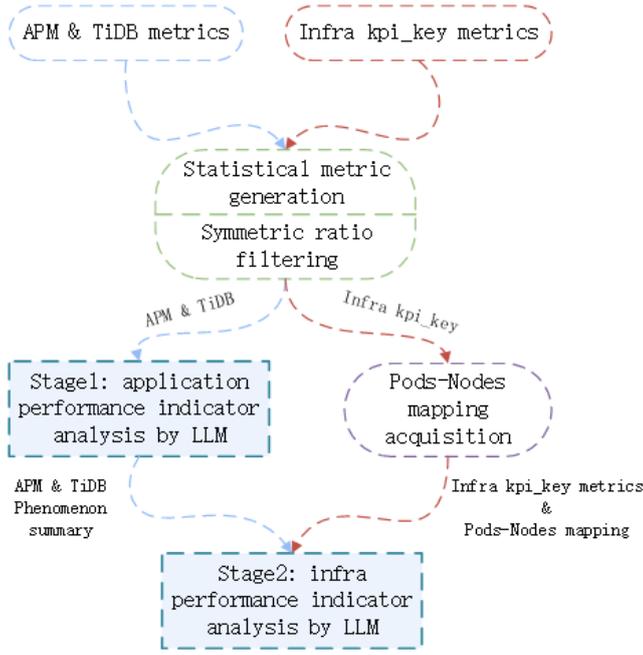


Figure 12: Schematic diagram of two-stage large language model phenomenon summary.

3.4.1 Principle and Method Construction of Large Language Model-Driven Phenomenon Summary

a Limitations of Traditional Anomaly Detection Methods

In distributed microservice architectures, monitoring metrics exhibit characteristics of high dimensionality, multi-level hierarchy, and strong correlations. Traditional anomaly detection methods, whether threshold-based or machine learning-based, often only identify numerical anomalies in individual metrics, struggling to capture the correlations between metrics and their business semantics. Furthermore, these methods lack the ability to generate semantic descriptions of anomalous phenomena, failing to provide interpretable phenomenon summaries for subsequent multimodal fault root cause localization.

b Phenomenon Summary Method Based on Large Language Models

In the metric data processing method of this solution, a large language model is employed as the core tool for phenomenon summary, fully leveraging its robust semantic understanding and reasoning-inductive capabilities. By inputting the filtered monitoring data, the model automatically identifies anomalous variation patterns, correlation relationships, and business impacts, and outputs phenomenon descriptions with rich semantic content. It is particularly emphasized that the large language model in this module is only used for phenomenon summary and pattern description, and does not perform fault judgment or root cause inference. This lays a foundation for subsequent comprehensive root cause analysis that integrates multimodal data such as logs and traces.

c Definition of Normal Time Period

The system defines the normal time period using a relative time window approach. The specific rules are as follows: the normal

operation time periods are extracted as two intervals: one from 10 minutes after the end of the previous fault to the start of the current fault, and the other from 10 minutes after the end of the current fault to the start of the next fault. This definition method not only ensures that the data in the normal time period reflects the stable operation state of the microservice system (avoiding interference from the "aftermath" of faults) but also effectively reduces the impact of business fluctuations during different time periods of the day (e.g., high visit volume during the daytime and low visit volume at night) by selecting adjacent time windows before and after fault periods. In turn, it provides a more reliable performance baseline for comparative analysis.

3.4.2 Multi-Level Monitoring Metric System and Filtering

a Pod-Service Unified Analysis Architecture

The proposed system uniformly processes the pod and service levels, and automatically extracts service identifiers by parsing pod names. For instance, pod instances such as frontend-0, frontend-1, and frontend-2 are uniformly classified under the "frontend" service by removing their numeric suffixes. This design simplifies the analysis complexity in microservice architectures while maintaining dual perspectives at both the service level and instance level, enabling effective identification of service-level common issues and specific anomalies in individual instances.

b Application Performance Monitoring Metric Filtering

From the numerous redundant metrics in microservice APM monitoring data, 7 core key metrics are carefully selected: client_error_ratio (client error rate), error_ratio (overall error rate), request (number of requests), response (number of responses), rrt (average response time), server_error_ratio (server error rate), and timeout (number of timeouts). These metrics cover the core dimensions of microservice performance, enabling comprehensive reflection of a service's health status and anomalous characteristics while avoiding the impact of redundant data on the effectiveness of phenomenon analysis.

c Hierarchical Coverage of Infra Performance Monitoring Metrics

An infrastructure monitoring system covering both container and node layers is constructed. The container layer includes 9 pod infrastructure metrics: pod_cpu_usage (pod CPU usage rate), pod_memory_working_set_bytes (pod memory working set size), pod_fs_reads_bytes (pod filesystem read bytes), pod_fs_writes_bytes (pod filesystem written bytes), pod_network_receive_bytes (pod network received bytes), pod_network_receive_packets (pod network received packets), pod_network_transmit_bytes (pod network transmitted bytes), pod_network_transmit_packets (pod network transmitted packets), and pod_processes (number of running processes in pod). The node layer comprises 16 node infrastructure metrics: node_cpu_usage_rate (node CPU usage rate), node_memory_usage_rate (node memory usage rate), node_disk_read_bytes_total (total disk read bytes), node_disk_written_bytes_total (total disk written bytes), node_disk_read_time_seconds_total (total disk read time in seconds), node_disk_write_time_seconds_total (total disk write time in seconds), node_filesystem_usage_rate (filesystem usage rate),

node_network_receive_bytes_total (total network received bytes), node_network_transmit_bytes_total (total network transmitted bytes), node_network_receive_packets_total (total network received packets), node_network_transmit_packets_total (total network transmitted packets), and node_sockstat_TCP_inuse (number of TCP connections), among others. Hierarchical monitoring ensures comprehensive coverage from microservice applications to the underlying infrastructure of specific deployments.

d Monitoring Metrics for TiDB Database Components

A specialized monitoring metric system is constructed for the TiDB distributed database, covering three core components: tidb-tidb, tidb-tikv, and tidb-pd. It includes database-specific metrics such as failed_query_ops (number of failed queries), duration_99th (99th percentile request latency), connection_count (number of connections), server_is_up (number of alive service nodes), cpu_usage (CPU usage rate), memory_usage (memory usage), store_up_count (number of healthy Stores), store_down_count (number of down Stores), store_unhealth_count (number of unhealthy Stores), storage_used_ratio (used capacity ratio), available_size (available storage capacity), raft_propose_wait (RaftPropose waiting latency), raft_apply_wait (RaftApply waiting latency), and rocksdb_write_stall (number of RocksDB write stalls). This enables professional monitoring and analysis of the data storage layer.

3.4.3 Data Processing and Anomaly Filtering Workflow

a Hierarchical Data Loading and Organization

The system constructs distinct data paths for file localization based on monitoring levels, using each input date information output from the preprocessing stage. Specifically, APM business monitoring data is stored in the directory "Year-Month-Day/metric-parquet/apm/pod", infrastructure node monitoring data resides in "Year-Month-Day/metric-parquet/infra/infra_node", container monitoring data is located in "Year-Month-Day/metric-parquet/infra/infra_pod", and TiDB database monitoring data is stored in directories such as "Year-Month-Day/metric-parquet/infra/infra_tidb". Through this hierarchical file organization method, precise data localization and loading are realized in accordance with monitoring dimensions.

b Precise Time Window Filtering and Data Merging

The system performs nanosecond-level time window filtering on monitoring data at all levels to ensure that the analysis scope fully covers the fault time period. Data from two normal time periods before and after the fault period are merged. The 10 minutes immediately preceding the fault period are excluded to eliminate the impact of fault "aftermath", and random fluctuation effects are removed by excluding the two maximum and two minimum extreme values, thereby constructing a stable statistical baseline. This processing method ensures the representativeness and stability of data during normal periods.

c Significance Filtering Based on Statistical Symmetric Ratio

A significance filtering mechanism based on statistical symmetric ratios is implemented, which calculates the symmetric ratios of medians and 99th percentiles between the fault period and normal

periods. Stable metrics with a variation amplitude of less than 5% are automatically filtered out, while only potentially critical abnormal metrics with significant changes are retained for the large language model analysis process. This filtering strategy can compress massive raw data into a set of high-value abnormal features, substantially reducing the context of the large language model by approximately 50% in terms of token count. Meanwhile, it significantly improves the efficiency and quality of phenomenon summarization. Taking the P50 symmetric ratio as an example, its specific formula is defined as follows:

$$P50_{\text{symmetric-ratio}} = \frac{|P50_{\text{fault}} - P50_{\text{normal}}|}{\frac{P50_{\text{fault}} + P50_{\text{normal}}}{2} + \varepsilon} \quad (1)$$

where $P50$ denotes the median, $P50_{\text{fault}}$ and $P50_{\text{normal}}$ represent the medians during the fault period and normal periods respectively, and ε is a minimal value.

For the filtered key metrics, the pandas 'describe' method is utilized to extract comprehensive descriptive statistical information, including key statistics such as median, interquartile range, and 99th percentile. Through a standardized data representation format, monitoring metrics of different types and dimensions are converted into structured inputs understandable to large language models, ensuring the accuracy and consistency of subsequent phenomenon summarization.

3.4.4 Two-Stage Large Language Model Phenomenon Summarization and Output Generation

a Stage 1: Application Performance Monitoring Phenomenon Identification and Summarization

The system first conducts a comprehensive analysis of microservice application performance monitoring metrics and TiDB database component metrics, organizing filtered key abnormal metrics by service type where each service includes data from its multiple subordinate Pod instances; for example, the emailservice contains comparative APM metrics such as client_error_ratio (client error rate), error_ratio (overall error rate), request (number of requests), response (number of responses), rrt (average response time), server_error_ratio (server error rate), and timeout (number of timeouts) across Pod instances like emailservice-0, emailservice-1, and emailservice-2, while integrating database-specific metrics of TiDB components including failed_query_ops (number of failed queries), duration_99th (99th percentile request latency), and connection_count (number of connections). A structured table with comparative data between normal periods and the fault period is constructed, and using a service-Pod hierarchical organization method, the large language model is invoked for the first-round phenomenon summarization (with prompt design shown in Figure 13); the data comparison table adopts JSON format, which, compared with Markdown, maintains clarity while significantly reducing token usage (specific format presented in Figure 14), and while analyzing Pod-level phenomena, the system naturally summarizes service-level common phenomena and differentiated performances, with the large language model only describing phenomena at this stage without making fault judgments.

<pre> 请根据提供的APM（应用性能监控）指标数据和TiDB分布式数据库指标数据，描述所有服务在正常期间和故障期间的业务服务性能表现差异现象。 ## 系统整体信息 - **微服务总数**：{microservice_count} - **TiDB组件数**：{tidb_service_count} - **服务总数**：{len(all_services)} - **Pod总数**：{total_pods} - **服务列表**：{all_services} ## APM关键指标说明（微服务） ### 请求响应类指标：{...} ### 异常类指标：{...} ## TiDB关键指标说明（数据库组件） ### TiDB组件指标：{...} ### TiKV组件指标：{...} ### PD组件指标：{...} ## 数据对比表格,特别注意**缺失数据和空数据,代表数据波动极小,通常情况默认正常**： {json.dumps(combined_json,ensure_ascii=False,indent=2)} </pre>	<pre> ## 现象描述要求 请从以下维度进行现象描述，**仅描述观察到的现象，不做异常判断或结论**： ### 微服务级别现象观察 - 集中在同一类型的微服务中出现的问题，比如emailservice-0, emailservice-1, emailservice-2, 都存在相似的异常数据变化 - 个别Pod的异常表现特征，比如cartservice-0中存在异常的数据变化，而cartservice-1, cartservice-2以及其他pod正常 ### TiDB数据库组件现象观察 - TiDB组件（tidb-tidb）、TiKV组件（tidb-tikv）和PD组件（tidb-pd）的异常数据变化 ## 重要提示 **这是为后期综合决策分析提供的系统级现象总结，请控制总结内容在2000字左右，重点突出主要变化现象。 ** **总结要求： ** - 如果整体表现正常稳定，请简要说明“系统各项指标表现稳定，未观察到显著变化现象” - 如果出现明显变化，请重点描述变化显著的服务、指标和现象 - 采用概括性语言，重点关注对业务影响较大的指标变化 - 缺失或为空的数据表示波动极小，可视为正常，无需描述 请基于APM业务监控数据和TiDB数据库监控数据客观描述观察到的现象，控制在2000字以内，为后续综合分析提供简洁有效的现象总结。 </pre>
--	--

Figure 13: Prompt for stage 1’s comprehensive phenomenon summarization (The prompt is designed in Chinese; see source code and modify as needed).

```

{
  服务1 (如adservice): {
    服务名称: 服务1,
    服务类型: microservice 或 TiDB-component,
    pod数量: 3,
    pod列表: [Pod-1, Pod-2, Pod-3],
    具体pods指标: {
      Pod-1 (如adservice-1): {
        具体指标 (如error_ratio): {正常期间中位数: xx.xx, 正常期间四分位距: xx.xx, 正常期间99分位数: xx.xx,
          故障期间中位数: xx.xx, 故障期间四分位距: xx.xx, 故障期间99分位数: xx.xx },
        其他指标: {...}, .....
      }
      Pod-2: {...},
      Pod-3: {...},
    },
    服务2: {...},
    ...
  }
}

```

Figure 14: Example of comparative data for business metrics and TiDB database components.

b Stage 2: Comprehensive Phenomenon Summarization and Correlation Analysis of Infrastructure Machine Performance Metrics

Building on the analysis of application performance monitoring data, the system further integrates machine performance monitoring data from the infrastructure layer, encompassing both the pod container layer and node layer, and by combining the deployment topology information of pods on each node with the service analysis results from Stage 1, it constructs a full-stack phenomenon view covering

three levels—application, container, and node—integrating container metrics such as pod_cpu_usage (Pod CPU usage rate), pod_memory_working_set_bytes (Pod memory working set size), and pod_network_receive_bytes (Pod network received bytes), as well as node metrics including node_cpu_usage_rate (node CPU usage rate), node_memory_usage_rate (node memory usage rate), node_disk_read_bytes_total (total disk read bytes), and node_network_transmit_bytes_total (total network transmitted bytes); the large language model is then invoked to conduct the second-round comprehensive phenomenon summarization, with data comparison tables also adopting JSON format that, after de-

```

{
metric的kpi_key指标数量: 15,
总metric的kpi_key列表: [...],
Nodes层次信息: {
  节点1 (如aiops-k8s-01):
  {
    节点名称: 节点1,
    节点上部署pod数量: 8,
    部署的pods具体名称: [...],
    该节点metrics的kpi_key具体统计信息:
    {
      指标1 (如node_cpu_usage_rate): {正常期间中位数: xx.xx, 正常期间四分位距: xx.xx, 正常期间99分位数: xx.xx,
      故障期间中位数: xx.xx, 故障期间四分位距: xx.xx, 故障期间99分位数: xx.xx },
      其他kpi_key指标: {...},
      ...
    },
    该节点上所部署的pods详细信息:
    {
      Pod-1 (如productcatalogservice-1):
      {
        指标1 (如pod_memory_working_set_bytes): {...},
        其他kpi_key指标: {...}
      },
      其他Pod: {...},
      ...
    },
    其他节点: {...},
    ...
  }
}

```

Figure 15: Example of statistical data comparison for kpi_key metrics between different nodes and the pods deployed on them.

<pre> 请基于提供的Service级别分析结果、Pod部署信息和基础设施监控指标数据，进行全局的现象总结分析。 ## Service级别分析结果回顾 {service_analysis_result} ## 集群基础设施信息 - **节点总数**：{len(all_nodes)} - **监控指标总数**：{len(all_metrics)} - **集群总Pod数**：{total_pods} - **节点列表**：{all_nodes} - **监控指标**：{[metric_chinese_names.get(m, m) for m in all_metrics]} ## 使用规范说明 • 所有监控指标均为 'kpi_key' 指标（如 'node_cpu_usage_rate'），请始终使用这些原始英文名称进行分析与输出； • 严禁使用中文或缩写形式（如“CPU”、“CPU使用率”、“磁盘读写”等）代替； • 必须显式包含对应的 'kpi_key'，每次提及 'kpi_key' 指标如('node_cpu_usage_rate') 时，必须显式包含对应的 'kpi_key'，必须显式指出这是一个 'kpi_key' 指标如： • 错误示例：节点级指标变化幅度（CPU ↑23%） • 正确示例：kpi_key指标'node_cpu_usage_rate' 在节点 aiops-k8s-08 上上升了 23% ## 基础设施指标分类说明 ### 计算资源类指标(kpi_key): {...} ### 存储资源类指标(kpi_key): {...} ### 网络资源类指标(kpi_key): {...} ## 基础设施指标数据对比表格,特别注意**缺失数据和空数据,代表数据波动极小,通常情况默认正常**： {json.dumps(combined_json, ensure_ascii=False, indent=2)} </pre>	<pre> ## 综合现象分析要求 请从以下维度进行全局现象描述，**仅描述观察到的现象，不做异常判断或结论**： ### 1. Node级别现象观察 基于Node的正常时间段与异常时间段的数据对比，描述： - 同一Node的正常时间段和异常时间段，显著的指标异常变化或显著的 'kpi_key' 指标异常变化 - 不同Node之间比较表现出的异常差异 ### 2. Service级别现象观察 - 集中在同一类型的服务中出现的问题，比如emailservice-0, emailservice-1, emailservice-2, 都存在相似的异常数据变化，则描述具体变化现象，因为这可能是emailservice存在潜在问题 ### 3. Pod级别现象观察 - 个别Pod的异常表现特征 - 如 cartservice-0 中存在异常的数据变化，而cartservice-1, cartservice-2以及其他pod正常，则可能是单独的pod级别的异常现象 - 大多数异常的Pod是否部署在同一个Node，还是分散在不同的Node，是否存在Node级别的异常现象 ## 重要提示 **这是为后期多模态综合决策分析提供的全局现象总结，请控制总结内容在2000字左右，重点突出主要变化现象。** **总结要求**： - 必须在输出中包含原始 'kpi_key' 指标名称，如 'node_cpu_usage_rate' - 分析原因时必须明确指出该原因属于哪个 'metric'（其他指标 "以及/或者" 哪个 'kpi_key'）下的观察 - 如果出现明显变化，请重点描述变化显著的服务、指标和现象 - 采用概括性语言，重点关注对业务影响较大的指标变化 - 关注异常现象倾向于node, service还是pod级别的异常 - 提供系统级的综合现象描述，为后续决策提供全面视角 - 采用客观、描述性的语言，避免主观判断 - 缺失或为空的数据表示波动极小，可视为正常，无需描述 请基于Service分析、Pod部署信息和基础设施监控数据，提供全局的综合现象总结，控制在2000字以内。 </pre>
---	---

Figure 16: Prompt for the stage 2's comprehensive phenomenon summarization.

<pre> #### 应用性能异常现象和基础设施机器性能异常现象综合分析 #### 1. Node级别现象观察 **计算资源类指标:** - kpi_key指标`node_cpu_usage_rate`在多个节点出现显著变化: - aiops-k8s-03: 从正常中位数35.86%升至39.97% (+11.5%), 99分位数从39.15%升至44.2% - aiops-k8s-05: 从14.95%升至21.62% (+44.6%), 99分位数从18.99%升至26.32% - aiops-k8s-07和aiops-k8s-08保持相对稳定(波动<±5%) - kpi_key指标`node_memory_usage_rate`呈现差异化变化: - aiops-k8s-01: 从50.26%升至54.94% (+9.3%), 伴随kpi_key指标`node_memory_MemAvailable_bytes`下降9.2% - aiops-k8s-08: 从55.92%升至65.2% (+16.6%), 可用内存下降22.6% - aiops-k8s-03和aiops-k8s-04出现内存使用率下降现象(分别-4.7%和-25.8%) **存储资源类指标:** - kpi_key指标`node_disk_written_bytes_total`在多个节点出现异常: - aiops-k8s-03: 中位数从60,620.8升至145,681.07 (+140%), 99分位数增长32.5% - aiops-k8s-08: 写入量保持高位(中位数>150MB), 但波动范围缩小 - aiops-k8s-02写入量异常下降50% **网络资源类指标:** - kpi_key指标`node_network_transmit_packets_total`在aiops-k8s-08出现基线偏移(中位数1.2→1.27) - 其他网络指标(如`node_network_receive_bytes_total`)未显示显著变化 **节点间差异:** - aiops-k8s-03和aiops-k8s-05同时出现CPU和磁盘写入压力 - aiops-k8s-01和aiops-k8s-08主要面临内存压力 - aiops-k8s-04和aiops-k8s-07资源指标相对稳定 #### 2. Service级别现象观察 **广告服务(adservice):** - adservice-0表现最显著: - RRT中位数增长6.5倍(3,230 μs→21,092 μs) - kpi_key指标`pod_cpu_usage`从0.01飙升至0.6(60倍增长) - kpi_key指标`pod_memory_working_set_bytes`增长27.1% - adservice-1和adservice-2未显示相同程度的CPU/内存压力 **购物车服务(cartservice):** - cartservice-0和cartservice-2: - RRT 99分位数分别增长4倍和43% - 但kpi_key指标`pod_memory_working_set_bytes`仅增长2.2%-4.9% - 未观察到对应节点的CPU/内存瓶颈 - cartservice-1表现正常 **前端服务(frontend):** - 全实例RRT增长(12,000 μs→15,800-17,200 μs) - frontend-0的kpi_key指标`pod_memory_working_set_bytes`波动增大(IQR从185,185→594,384) - 部署在aiops-k8s-03/04/07的实例表现一致 </pre>	<pre> **缓存服务(redis-cart-0):** - timeout次数中位数从1次增至4次 - kpi_key指标`pod_memory_working_set_bytes`异常激增(中位数124.33→42,411.29) - kpi_key指标`pod_processes`的99分位数从1.52升至2.0 #### 3. Pod级别现象观察 **异常Pod分布特征:** - 关键异常Pod集中在特定节点: - aiops-k8s-03: adservice-0, frontend-0 - aiops-k8s-08: redis-cart-0, cartservice-2 - aiops-k8s-07: emailservice-0(RRT增长2.5倍) **典型异常Pod:** - adservice-0(aiops-k8s-03): - 唯一显示CPU使用率飙升的Pod(kpi_key指标`pod_cpu_usage` 0.01→0.6) - 内存使用与磁盘写入(kpi_key指标`pod_fs_writes_bytes`)同步增长 - redis-cart-0(aiops-k8s-08): - 内存占用99分位数从3.3GB升至5.3GB - 进程数(kpi_key指标`pod_processes`)达到上限阈值 - emailservice-0(aiops-k8s-07): - 网络吞吐量(kpi_key指标`pod_network_transmit_bytes`)增长35.3% - 但节点级网络指标无对应变化 **稳定服务特征:** - payment-service所有实例: - RRT波动<±10% - 内存使用(kpi_key指标`pod_memory_working_set_bytes`)无显著变化 - recommendation-service: - 所有实例指标波动在正常范围内 - 网络流量(kpi_key指标`pod_network_receive_packets`)变化<±5% #### 跨层级关联现象模式 1. **资源热点集中:** - CPU压力集中在aiops-k8s-03和aiops-k8s-05 - 内存压力集中在aiops-k8s-01和aiops-k8s-08 - 磁盘写入热点出现在aiops-k8s-03和aiops-k8s-08 2. **服务级联影响:** - 前端服务实例延迟增长可能由下游adservice/cartservice异常放大 - redis-cart-0的异常直接影响所有cartservice实例 3. **异常分布模式:** - 有状态服务(adservice/cartservice)问题更显著 - 同一服务的不同实例表现差异大(如adservice-0 vs adservice-1) - 无状态服务(payment-service等)整体稳定 4. **指标关联现象:** - adservice-0的CPU飙升与RRT增长直接相关 - redis-cart-0内存增长与timeout次数增加同步出现 - 节点级磁盘写入增长(如aiops-k8s-03)未对应到具体Pod的写入指标异常 </pre>
---	--

Figure 17: Output example of the metric fault summarization module.

describing the statistical indicators of node-level kpi_key, proceeds to describe the kpi_key indicators related to pods deployed on that node to ensure a coherent data flow (specific format shown in Figure 15 and prompt design presented in Figure 16), emphasizing the analysis of cross-level abnormal correlation relationships to output a complete fault phenomenon description spanning from microservice applications to infrastructure, while this stage also only performs phenomenon summarization to lay the foundation for subsequent multimodal root cause analysis.

c Structured Phenomenon Output

After processing via the two-stage large language model analysis, the metric fault summarization module outputs a phenomenon description that includes the following content:

Application Performance Anomaly Phenomena: Leveraging the service-Pod hierarchical information, microservice components and database components with significant performance changes are identified. This includes changes in request processing metrics such as request (number of requests), response (number of responses), and rrt (average response time); fluctuations in abnormal metrics like error_ratio (overall error rate), client_error_ratio (client error rate), server_error_ratio (server error rate), and timeout (number of timeouts); and variations in database performance metrics of TiDB components, including failed_query_ops (number of failed queries), duration_99th (99th percentile request latency),

and connection_count (number of connections). This enables the simultaneous representation of overall trends at the service level and individual differences at the Pod level.

Infrastructure Machine Performance Anomaly Phenomena: This section describes changes in resource status at the container and node levels, including variation patterns of computing resource metrics such as node_cpu_usage_rate (node CPU usage rate) and pod_cpu_usage (pod CPU usage rate); memory resource metrics like node_memory_usage_rate (node memory usage rate) and pod_memory_working_set_bytes (pod memory working set size); storage resource metrics including node_disk_read_bytes_total (total disk read bytes), node_disk_written_bytes_total (total disk written bytes), and pod_fs_reads_bytes (pod file system read bytes); and network resource metrics such as node_network_receive_bytes_total (total network received bytes), node_network_transmit_bytes_total (total network transmitted bytes), and pod_network_receive_bytes (pod network received bytes), as well as spatial characteristics of anomalies distributed across different nodes and containers.

Cross-Level Correlated Phenomenon Patterns: Through Pod deployment topology and phenomenon correlation analysis, the correlation patterns between service anomalies and infrastructure anomalies are identified, and the distribution characteristics and propagation paths of abnormal phenomena are determined. This provides critical phenomenon clues for subsequent multimodal fault root cause analysis that integrates log and trace data.

<pre> ### Language Enforcement -Input may contain Chinese. **but output MUST be entirely in English** (no Chinese characters). 请根据提供的 {modalities_text}, 进行综合故障分析, 识别最有可能的单一故障原因。不要包含任何其他解释或文本。 特别注意**缺失数据和空数据,代表数据波动极小,通常情况默认正常,严禁分析和定位为根因** ### 微服务架构调用关系图谱 理解以下关键调用路径有助于识别故障传播和根因定位: **主要调用路径:** 1. **用户请求入口:** User → frontend (所有用户请求的统一入口) 2. **购物核心流程:** frontend → checkoutservice → (paymentservice, emailservice, shippingservice, currencyservice) 3. **商品浏览相关:** frontend → (adservice, recommendationservice, productcatalogservice, cartservice) 4. **服务间依赖:** recommendationservice → productcatalogservice (推荐依赖商品目录) 5. **数据存储层:** - adservice/productcatalogservice → tidb (广告和商品数据存储) - cartservice → redis-cart (购物车缓存) - tidb 集群内部: tidb → (tidb-tidb, tidb-tikv, tidb-pd) **故障传播分析要点:** - **上游故障影响:** frontend故障会影响所有下游服务 - **checkout相关:** checkoutservice故障会影响支付、邮件、物流、货币服务 - **数据层故障:** tidb故障影响adservice和productcatalogservice; redis故障影响cartservice - **横向依赖:** recommendationservice依赖productcatalogservice 要求: ** 1. 综合多种监控数据进行分析, 优先考虑数据间的关联性, **特别关注调用路径上的故障传播模式** 2. 只返回一个最可能的故障分析结果 3. 故障级别判断标准: **Node级别故障:** 单个节点的监控指标 (kpi_key) (node_cpu_usage_rate,node_filesystem_usage_rate等) 对比正常期间,故障期间存在显著异常变化, 且该节点上的多个不同服务的Pod均受影响 **Service级别故障:** 同一服务的多个Pod实例 (如emailservice-0, emailservice-1, emailservice-2) 都出现相似的异常数据变化, 表明服务本身存在问题 **Pod级别故障:** 单个Pod (如cartservice-0) 出现异常数据变化, 而同服务的其他 Pod (cartservice-1, cartservice-2) 及其他Pod正常 **重要说明:** 所有监控指标均为 `kpi_key` 指标 (例如 `node_cpu_usage_rate`), 请在描述中直接使用这些原始 `kpi_key` 英文指标名, 不得使用中文或其他名称。 </pre>	<pre> 4. 请确保: - component必须从提供的组件列表中选择, 组件列表包含三种故障层级: * 节点名(aiops-k8s-01-08) - 表示节点级别的基础设施故障 * 服务名(cartservice等) - 表示微服务级别的故障 * Pod名(cartservice-0等) - 表示单个Pod级别的故障 - reason必须说明在哪个具体的监控指标(kpi_key)发生故障,且要使用kpi_key原始英文 (node_cpu_usage_rate,node_filesystem_usage_rate等),要简明扼要, 说明故障的根本原因 - time要基于数据中的earliest_time - observation要基于多模态数据中的关键证据,要简明扼要,并明确指出来自哪个模态(如来自metric必须指出相应的监控指标(kpi_key),log只需要提及log(日志)检索和日志故障关键词,trace(调用链)需要按照*故障根因component*在*trace路径*出现*异常调用行为(caller/callee/self-loop)*的方式提及) - **重要:** 在分析trace数据时, 请参考上述调用关系图谱, 优先分析关键调用路径上的异常(如frontend→checkout→payment链路, 或数据存储相关的服务调用) - **特别要求:** 严禁分析和定位缺失数据和空数据为根因,默认其正常 5. The JSON output must be fully in English. Any Chinese characters are strictly prohibited. **Strictly follow the JSON format below**:</pre> <pre> { "component": "Select from the following components: {components_list}", "reason": "Most likely root cause based on comprehensive multi-modal analysis; (must include kpi_key for metrics. (Do not infer from missing data.))", "time": "Recommended time format: 'YYYY-MM-DD HH:mm:ss', e.g. '2025-04-21 12:18:00'", "observation": "Key evidence from multiple data sources (log,metric,trace)", "reasoning_trace": [{ "step": 1, "action": "Such as: LoadMetrics(checkoutservice)", "observation": "Describe (≤20 words) the most critical anomaly in metric modality, must include exact kpi_key and change (e.g., 'node_cpu_usage_rate' increased 35% at 12:18 in metric)" }, { "step": 2, "action": "Such as: TraceAnalysis('frontend -> checkoutservice -> paymentservice')", "observation": "Describe (≤20 words) the most critical abnormal behavior in trace modality, include trace path and anomaly type based on call graph (e.g., 'timeout in 'checkoutservice -> paymentservice' call in trace)" }, { "step": 3, "action": "Such as: LogSearch(checkoutservice)", "observation": "Describe (≤20 words) the most critical anomaly in log modality, mention error keyword and count/context (e.g., 'IOError found in 3 entries in log')" }] } </pre> <p>可用的监控数据: {data_content}</p>
--	---

Figure 18: Prompt of multimodal analysis.

Through the systematic processing workflow outlined above, the metric fault summarization module achieves the conversion of massive multi-dimensional monitoring data into concise phenomenon descriptions. By fully leveraging the semantic understanding and reasoning capabilities of large language models, this module focuses on phenomenon summarization rather than fault judgment. It not only significantly improves the accuracy and completeness of abnormal phenomenon identification but also provides a high-quality semanticized metric data foundation for subsequent fault root cause analysis that integrates multimodal data.

The output example of the metric fault summarization module is specifically shown in Figure 17.

3.5 Multimodal Root Cause Analysis

The multimodal root cause analysis module serves as the core decision-making component for fault root cause localization in mi-

croservice systems. It adopts a comprehensive reasoning strategy based on large language models, conducting integrated analysis of preprocessed fault log data, abnormal trace patterns, and summarized metric phenomena to achieve automated reasoning from multi-dimensional abnormal information to accurate root cause localization. By fully leveraging the cross-modal understanding and logical reasoning capabilities of large language models, this module provides final component localization and cause explanation for complex microservice faults.

3.5.1 Input Data Integration

The system first integrates the processing results of the three modalities. The log fault extraction module, based on the Drain log parsing algorithm combined with a multi-level data filtering mechanism, compresses massive raw logs into high-quality fault features. The trace fault detection module outputs two types of abnormal

```

{
  "fault_category": "jvm fault",
  "fault_type": "jvm cpu",
  "instance_type": "pod",
  "service": "adservice",
  "instance": ["adservice-0"],
  "start_time": "2025-06-12T08:02:48Z",
  "end_time": "2025-06-12T08:25:48Z",
  "uuid": "8c1e8ce9-237",
  "key_observations": [
    {
      "type": "log",
      "keyword": ["adservice--stress"], ["deadlineexceeded"], ["invocationtargetexception"], ["transformlistener"]
    },
    {
      "type": "trace",
      "subtype": "latency\_anomalies",
      "keyword": ["frontend"]
    },
    {
      "type": "metric",
      "keyword": ["pod\_cpu\_usage"], ["rrt", "rrt_max"]
    }
  ],
  "key_metrics": ["pod\_cpu\_usage"], ["adservice--stress"],
  "fault_description": ["JVM CPU spike", "JVM CPU overload", "JVM high CPU usage", "JVM CPU saturation"]
}

```

Figure 19: Good case - ground truth range.

```

{
  "component": "adservice-0",
  "uuid": "8c1e8ce9-237",
  "reason": "kpi_key `pod\_cpu\_usage` surged and kpi_key `pod_memory_working_set_bytes` increased significantly in adservice-0",
  "reasoning_trace": [
    {
      "step": 1,
      "action": "LoadMetrics(adservice-0)",
      "observation": "kpi_key `pod\_cpu\_usage` surged and kpi_key `pod_memory_working_set_bytes` increased in metric"
    },
    {
      "step": 2,
      "action": "TraceAnalysis('frontend -> adservice')",
      "observation": "hipstershop.AdService/GetAds duration spiked in trace"
    },
    {
      "step": 3,
      "action": "LogSearch(adservice-0)",
      "observation": "DeadlineExceeded error found in log"
    }
  ]
}

```

Figure 20: Good case - actual output values.

information: the duration anomaly detection results include the top 20 most frequent anomalies identified by Isolation Forest, covering call relationships, performance comparisons, and anomaly statistics; the status anomaly detection results include the top 20 most frequent status code anomalies, covering error types and specific descriptions. The metric fault summarization module outputs a phenomenon description summary analyzed by the two-stage large language model, which contains approximately 2,000 words of detailed abnormal phenomenon descriptions, covering monitoring anomaly patterns across node, service, and pod levels.

3.5.2 Multimodal Prompt Generation

Based on the integrated multimodal data, the system constructs a dedicated cross-modal analysis prompt template. The prompt design adheres to a structured principle, clearly identifying the source, type, and analysis focus of data from each modality to provide clear data understanding guidance for the large language model. Meanwhile, the output requirements and format specifications are explicitly defined in the prompt, ensuring the large language model can conduct systematic root cause reasoning based on the multimodal evidence chain. The prompt design is shown in Figure 18.

3.5.3 Structured Output Format Design

The system primarily guides the standardized output format by providing output example prompts, ensuring the consistency and interpretability of root cause analysis results. The output result includes three core fields: Fault Component (component), Fault Reason (reason), and Reasoning Process (reasoning_trace). The Fault Component field clearly identifies the name of the problematic microservice component; the Fault Reason field provides a concise and clear root cause description; and the Reasoning Process field details the model’s analytical logic and evidence chain, laying the foundation for result verification and subsequent optimization.

3.5.4 Result Extraction and Verification Mechanism

Due to the potential instability of large language model outputs (such as missing closing brackets, inclusion of redundant explanations, etc.), the system employs regular expressions and structured parsing methods to extract standardized results from the natural language outputs of large language models. It ensures the format correctness of output results through JSON format verification and guarantees the validity of key fields via content integrity checks. For results with parsing failures or format inconsistencies, the system provides retry and exception handling procedures to maximize the quality of outputs.

Through the multimodal root cause analysis workflow outlined above, the system achieves end-to-end automated processing from heterogeneous monitoring data to fault localization. By fully leveraging the advantages of large language models in cross-modal understanding and logical reasoning, this module provides an efficient, accurate, and interpretable solution for fault root cause localization in complex microservice systems.

4 Result Analysis

4.1 Comparison Between Good Cases and Bad Cases

As shown in Figure 19 and Figure 20, the proposed intelligent fault root cause localization system with multimodal data fusion can accurately identify the most likely faulty components and key root cause derivation evidence from the three modalities of monitoring data (logs, trace data, and system metrics) through in-depth integration. This demonstrates the effective implementation of the designed multimodal extraction scheme. Furthermore, the scheme is capable of making reasonable sequential reasoning, effectively identifying fault hierarchy types, and determining the actual hierarchical sources (pods, services, nodes) of potentially abnormal data indicators. Ultimately, the multimodal root cause analysis module can make reasonable fault root cause judgments based on log fault data, trace anomaly patterns, and metric phenomena. By effectively combining the multimodal analysis and logical reasoning capabilities of large language models, it achieves the goal of reasonable fault root cause localization.

Furthermore, we also explored typical bad cases to further identify areas for improvement in the proposed scheme. As shown in Figure 21 and Figure 22, the large language model exhibited hallucination during the final root cause analysis (indicated in blue in the figures). The relevant trace call chains were not actually input into the large language model for root cause analysis, leading to an incorrect reasoning chain. This may result in deviations in the final fault root cause judgment and affect the stability of the scheme. In subsequent work, optimizations may be implemented by adding

retrieval-augmented generation, jury mechanisms, and other methods. Additionally, after the competition released the answers, the log-related module can generate more effective keyword filtering templates through approaches such as enhancing the large language model’s semantic judgment using broader and more representative fault data, thereby improving log filtering capabilities.

4.2 Ablation Study

Table 1: Ablation Study of Each Component.

Log	Trace	Metric	Score
✓	×	×	23.59
×	✓	×	31.09
×	×	✓	42.78
✓	✓	×	35.32
×	✓	✓	48.58
✓	×	✓	51.27
✓	✓	✓	50.71

To further investigate the collaborative mechanism among the components of the proposed scheme, we conducted systematic ablation experiments. By comparing the performance of different modal combinations, we revealed the unique value and complementary relationships of each module. As shown in the experimental results in Table 1, in the single-modal experiments, the metric module performed the most prominently (42.78 points), mainly due to its two-stage LLM analysis strategy that effectively covers full-stack monitoring perspectives across the service, pod, and node levels; the trace module ranked second (31.09 points), demonstrating strong capabilities in call chain analysis and identifying dependencies between services; while the log module exhibited relatively weaker performance when used independently (23.59 points), primarily because the unstructured nature of log data limits the accuracy of its standalone analysis. The dual-modal combination experiments revealed significant synergistic effects: the log+metric combination achieved the best performance (51.27 points), with this strong synergy stemming from the fact that metrics provide quantitative localization of performance anomalies while logs offer rich semantic descriptions of faults—together forming a perfect complementarity from numerical anomalies to semantic explanations; the trace+metric combination ranked next (48.58 points) as they share a natural correlation in the performance monitoring dimension; and the log+trace combination showed relatively weaker performance (35.32 points) due to the lack of quantitative support from the metric module. The complete three-modal fusion system achieved 50.71 points, which, although slightly lower than the 51.27 points of the optimal dual-modal combination (log+metric), still significantly outperformed other combinations—indicating that the synergy between logs and metrics can cover most fault patterns, while the trace module retains unique value in scenarios involving complex call chain anomalies. Ultimately, the results of the ablation experiments fully validate the scientific rigor of the multimodal fusion architecture in this scheme, confirm the key role of the metric module as the system’s core analysis engine, demonstrate the optimal performance of the log+metric combi-

```

{
  "fault_category": "pod fault",
  "fault_type": "pod failure",
  "instance_type": "service",
  "service": "cartservice",
  "instance": "cartservice",
  "start_time": "2025-06-05T18:10:05Z",
  "end_time": "2025-06-05T18:34:05Z",
  "uuid": "38ee3d45-82",
  "key_observations": [
    {
      "type": "log",
      "keyword": ["Connection", "refused"], ["Error", "while", "dialing"], ["failure"], ["unavailable"]
    },
    {
      "type": "metric",
      "keyword": ["client_error", "client_error_ratio", "error", "error_ratio"]; ["pod_cpu_usage"];
        ["pod_network_receive_bytes", "pod_network_receive_packets",
          "pod_network_transmit_bytes", "pod_network_transmit_packets"]
    },
  ],
  "key_metrics": ["client_error", "client_error_ratio", "error", "error_ratio", "pod_processes", "request", "response"],
  "fault_description": ["pod crash", "pod terminated unexpectedly", "container exit", "pod failure"]
}

```

Figure 21: Bad case - ground truth range.

```

{
  "component": "cartservice",
  "uuid": "38ee3d45-82",
  "reason": "All cartservice instances show `pod_network_receive_bytes` dropped to zero and high client error rates, indicating complete network connectivity failure",
  "reasoning_trace": [
    {
      "step": 1,
      "action": "LoadMetrics(cartservice)",
      "observation": "`pod_network_receive_bytes` dropped to zero for all cartservice instances in metric"
    },
    {
      "step": 2,
      "action": "TraceAnalysis('frontend -> cartservice')",
      "observation": "Connection refused errors in all frontend to cartservice calls in trace"
    },
    {
      "step": 3,
      "action": "LogSearch(cartservice)",
      "observation": "Connection refused errors found in frontend logs when calling cartservice"
    }
  ]
}

```

Figure 22: Bad case - actual output values.

nation, and provide important empirical support for constructing an accurate and interpretable solution for microservice fault root cause localization based on large-model agents.

5 Conclusion

This solution constructs a large-model agent analysis framework for fault root cause localization in microservice systems, which con-

sists of five core modules: data preprocessing, log fault extraction, call chain anomaly detection, metric phenomenon summarization, and multimodal root cause analysis. First, the data preprocessing module realizes structured parsing of fault input and timestamp standardization, laying a foundation for multi-source data fusion. The log module leverages the Drain algorithm and a multi-level filtering mechanism to extract error templates and screen valid

structured log features. The call chain module combines the Isolation Forest model with status code checks to perform dual detection of duration and status anomalies. The metric module filters core business and infrastructure metrics, and uses large language models to conduct two-level phenomenon summarization, generating highly readable anomaly descriptions. Finally, the multimodal root cause analysis module fuses the outputs of log, trace, and metric data, and drives large language models to conduct comprehensive reasoning through well-designed prompts, outputting structured root cause analysis results that include faulty components, causes, and reasoning trace. Ablation experiments fully verify the complementary value of multimodal data and the effectiveness of the system architecture. This framework performs excellently in complex microservice fault scenarios, ultimately achieving a score of 50.71.

References

- Pinjia He, Jieming Zhu, Zibin Zheng, and Michael R Lyu. 2017. Drain: An online log parsing approach with fixed depth tree. In *2017 IEEE international conference on web services (ICWS)*, pages 33–40. IEEE.
- Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation forest. In *2008 eighth ieee international conference on data mining*, pages 413–422. IEEE.