

# LogEval: A comprehensive benchmark suite for LLMs in log analysis

Tianyu Cui<sup>1</sup> · Shiyu Ma<sup>1</sup> · Ziang Chen<sup>1</sup> · Tong Xiao<sup>2</sup> · Chenyu Zhao<sup>1</sup> · Shimin Tao<sup>3</sup> · Yilun Liu<sup>3</sup> · Shenglin Zhang<sup>1,4</sup> · Duoming Lin<sup>1</sup> · Changchang Liu<sup>1</sup> · Yuzhe Cai<sup>1</sup> · Weibin Meng<sup>3</sup> · Yongqian Sun<sup>1,5</sup> · Dan Pei<sup>2</sup>

Accepted: 7 July 2025

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2025

#### Abstract

Log analysis is vital in Artificial Intelligence for IT Operations (AIOps) and plays a crucial role in ensuring software reliability and system stability. However, challenges such as the absence of comprehensive evaluation standards, inconsistencies in benchmarking practices, and limited exploration of Large Language Models (LLMs) in log-related tasks persist. To address these issues, we introduce LogEval, a comprehensive benchmark designed to systematically evaluate LLMs' performance across four key log analysis tasks: log parsing, log anomaly detection, log fault diagnosis, and log summarization. LogEval systematically tackles these challenges through the following aspects: (i) it incorporates 4,000 publicly available log entries, spanning diverse tasks and providing a strong foundation for evaluating LLM performance; (ii) it utilizes standardized prompts in both English and Chinese to ensure consistent and objective evaluations, this benchmark covers two experimental paradigms: Naive question-answering (Q&A) and self-consistency (SC) Q&A, under both zero-shot and few-shot settings, while also considering inference efficiency and average token usage; (iii) it features an open-source, continuously updated platform (https://nkcs.iops.ai/LogEval/) that integrates new LLMs and user-uploaded production data, fostering reproducibility and adaptability in performance comparisons. The experimental results provide valuable insights into the varying strengths of LLMs across different tasks, highlighting opportunities for further optimization and innovation for LLMs in log analysis. Our code repository is available at https://github.com/LinDuoming/LogEval.

**Keywords** Log analysis · Benchmark suite · Large language models · Prompt engineering

Communicated by: Markus Borg.

Published online: 10 October 2025

Extended author information available on the last page of the article



#### 1 Introduction

With the rapid advancement of information technology, software systems have become essential to the operations of businesses and organizations (Cito et al. 2015). These systems generate vast amounts of log data that capture operational behavior, status changes, and potential failures (Li et al. 2020). As software systems grow in both scale and complexity, the manual log analysis performed by experts is becoming increasingly difficult and errorprone (Zhang et al. 2021; Nedelkoski et al. 2020; Wang et al. 2024; Zhong et al. 2024). It is not only time-consuming but also inefficient, often leading to delayed responses to critical system failures (Locke et al. 2022; Ma et al. 2024; Lin et al. 2016). Consequently, there is an urgent need for automated log analysis tools capable of quickly providing insights to ensure the reliability and stability of modern large-scale systems (He et al. 2021).

To meet these needs, various automated log analysis tools have been developed, focusing on four primary tasks: log parsing (Meng et al. 2020; Zhu et al. 2019; Liu et al. 2022; Coustié et al. 2020; Le and Zhang 2023; Xiao et al. 2020; Zhu et al. 2019; Wang et al. 2022), log anomaly detection (Du et al. 2017; Meng et al. 2019; Guo et al. 2021; Le and Zhang 2022; Zhao et al. 2021; Zhang et al. 2019; Du et al. 2021), log fault diagnosis (Jia et al. 2021; Zhou et al. 2019; Zhang et al. 2021; He et al. 2018; Liu et al. 2022; Ma et al. 2022; Luo et al. 2021; Zhou et al. 2020), and log summarization (Meng et al. 2023; Locke et al. 2022). In recent years, deep learning (DL) methods have been widely applied in log analysis to address the limitations of traditional approaches (He et al. 2021; Sui et al. 2023). Unlike traditional machine learning (ML) methods, deep learning is more effective at handling large-scale and complex log data (Ma et al. 2024; Le and Zhang 2022; Zhang et al. 2019), and can automatically extract features in an end-to-end manner, avoiding the constraints of manual feature engineering and fixed rule-based methods (Liu et al. 2019). While deep learning-based approaches offer significant improvements over traditional ML methods, they also come with certain challenges (Le and Zhang 2024). One major limitation is that deep learning models, especially PLMs, require substantial computational resources and large datasets for both pre-training and fine-tuning (Ma et al. 2024). Additionally, these models may struggle with domain-specific terminologies, such as the abbreviations commonly found in logs, log events of the same type exhibit different semantics and different log events share many similar words but exhibit opposite (He et al. 2024), which are not typically present in general language corpora. Furthermore, the variability of logs across different systems poses another challenge, as DL models often need to be retrained or fine-tuned frequently to effectively handle new log types.

To address the generalization challenge, data requirements, and retraining issues of DL-based methods, researchers have begun to explore the use of Large Language Models (LLMs) in log analysis tasks (Liu et al. 2024; Zhong et al. 2024), as LLMs have demonstrated outstanding performance in various natural language processing (NLP) tasks such as text generation, language translation, and sentiment analysis. LLMs like GPT-4 (OpenAI et al. 2024), LLaMA-2 2023, ChatGLM-4 (THUDM 2024), and Qwen-1.5 (Bai et al. 2023) have shown promising performance in these tasks. Nevertheless, with the diversification of LLMs, their performance varies across different tasks. As a result, selecting the most appropriate LLM for a given log analysis task has become an important consideration in both research and practice. However, in the field of log analysis, there is currently no comprehensive and systematic evaluation standard or toolkit to help researchers and developers understand and compare the performance of different LLMs across various log analysis



tasks. The diverse architectures, model sizes, and applicability of LLMs make selecting the optimal model a complex task that lacks scientific guidance. Therefore, there is an urgent need to construct a unified benchmark that can scientifically assess the performance of different LLMs and provide objective, comprehensive comparisons. To address this, we propose and develop a comprehensive benchmark suite called *LogEval*, designed to evaluate LLMs' performance across various log analysis tasks. The main contributions of this paper are as follows:

- Diverse Log Dataset Collection: LogEval incorporates log datasets from multiple sources, covering core tasks such as log parsing, log anomaly detection, log fault diagnosis, and log summarization. This curated dataset provides a robust foundation for evaluating LLM performance comprehensively.
- Unified and Reliable Evaluation: LogEval utilizes standardized English and Chinese
  prompts to ensure consistent and objective assessments of LLMs. A unified prompt design is introduced for all tasks, minimizing variations caused by differing prompt styles
  and ensuring fair comparisons. The evaluation spans two paradigms: Naive Q&A and
  Self-Consistency Q&A under zero-shot setting and few-shot setting, while also considering inference efficiency and token usage.
- Dynamic Benchmarking Platform: LogEval features an open-source, continuously updated online platform (https://nkcs.iops.ai/LogEval/) that allows dynamic integration of new LLMs and user-uploaded production log data. This platform promotes reproducibility, fairness, and adaptability in performance comparisons, ensuring long-term relevance and scalability. Our code repository is available at https://github.com/LinDuoming/LogEval.

# 2 Background

Automated log analysis typically involves four core tasks—log parsing, log anomaly detection, log fault diagnosis, and log summarization. Each task addresses specific challenges and plays a crucial role in transforming raw log sequences into actionable insights. Figure 1 illustrates the sequence of these tasks in the log analysis pipeline. Below, we describe each task, the results depicted in the figure, and the evaluation metrics used to measure their performance.

Log Parsing Log parsing is the initial task in the log analysis pipeline. It involves transforming raw log data into a structured format that can be processed by subsequent tasks. The goal of log parsing is to extract relevant components (such as interface states, error messages, or system events) from unstructured logs and represent them in a consistent format, often as key-value pairs or predefined templates. In Fig. 1, the first section illustrates log entries such as interface state changes or member port status updates, which are parsed and organized

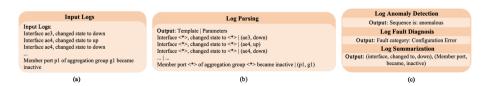


Fig. 1 A demonstration of the log analysis tasks



into templates (e.g., "Interface <\*>, changed state to <\*>") in the following table. This structured output enables the identification of recurring patterns or anomalies. **Metrics and Formula**: The parsing performance is evaluated using two key metrics:

• Parsing Accuracy: The formula for parsing accuracy is given by

$$Accuracy = \frac{Correctly Parsed Entries}{Total Entries} \times 100\%$$
 (1)

where **Correctly Parsed Entries** refers to log lines that exactly match predefined templates, and **Total Entries** refers to all log entries in the dataset.

• Edit Distance: The formula for edit distance is given by

Edit Distance = 
$$I + D + S$$
 (2)

where **Insertions** (**I**) is the number of characters added to match the template, **Deletions** (**D**) is the number of characters removed to match the template, and **Substitutions** (**S**) is the number of character replacements needed.

Log Anomaly Detection Once logs are parsed, the next step is log anomaly detection. This task aims to identify unusual log entries that may indicate potential issues or system faults. Anomalies are detected based on patterns that deviate from normal system behavior, such as unexpected state changes, errors, or performance issues. In Fig. 1, the second section of the diagram shows the log anomaly detection task. After parsing, each log sequence is examined for anomalous behavior. For example, an unexpected interface state change or an error message could be flagged as anomalous. The detected anomalies are then passed on for further investigation in the fault diagnosis stage. Metrics and Formula: The performance of anomaly detection is commonly assessed using the following metrics: - Precision (the proportion of detected anomalies that are true positives):

$$Precision = \frac{True \ Positives}{True \ Positives + False \ Positives}$$
(3)

Recall (the proportion of actual anomalies that are correctly detected):

$$Recall = \frac{True Positives}{True Positives + False Negatives}$$
 (4)

F1-score (the harmonic mean of precision and recall):

$$F1-score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$
 (5)



Where True Positives (TP) are correctly identified anomalies, False Positives (FP) are incorrectly flagged as anomalies, and False Negatives (FN) are missed anomalies.

Log Fault Diagnosis The log fault diagnosis task aims to identify the root cause of the detected anomalies by analyzing the correlations between log entries. This step often involves mapping anomalies to known fault categories or failure modes. In Fig. 1, the third section illustrates this task, where the system correlates parsed and anomalous log entries to diagnose faults. For example, a change in interface state from "up" to "down" might correlate with hardware failure or misconfiguration. Metrics: This task shares the same evaluation metrics with anomaly detection.

Log Summarization Finally, log summarization condenses large volumes of log data into a concise and interpretable format, highlighting key events that require attention. The goal is to present a summary of the most critical log entries in a format that is easy for system administrators to understand and act upon. In Fig. 1, the log summarization task is shown in the final section, where the relevant log entries identified during the previous stages are summarized. For example, entries like "interface changed state to down" and "member port became inactive" are distilled into a more concise format that provides key insights for further analysis. Metrics and Formula: The performance of summarization is commonly assessed using the following metrics:

• ROUGE-L F1: The formula for ROUGE-L is given by

$$ROUGE-L = \frac{|LCS(S, R)|}{|R|} \quad (Recall)$$
 (6)

and the formula for F1 is given by

$$F1 = 2 \times \frac{\text{ROUGE-L} \times P_{\text{LCS}}}{\text{ROUGE-L} + P_{\text{LCS}}}$$
 (7)

where **LCS** refers to the Longest Common Subsequence between summary (S) and reference (R), and **P\_LCS** refers to the Precision of LCS, which is  $\frac{|LCS(S,R)|}{|S|}$ .

• Threshold Accuracy: The formula for threshold accuracy is given by

Accuracy = 
$$\frac{\sum_{i=1}^{n} \mathbb{I}(\text{ROUGE-L}_i \ge \theta)}{n} \times 100\%$$
 (8)

where  $\theta$  is the similarity threshold (typically 0.7-0.9), and  $\mathbb{I}$  is the indicator function (1 if condition met, 0 otherwise).

In this study, we carefully selected evaluation metrics tailored to the nature of each log analysis task to ensure a comprehensive, objective, and practically meaningful assessment. For tasks such as log anomaly detection and log fault diagnosis, we adopt Accuracy and F1-score as the primary evaluation metrics. Accuracy reflects the overall correctness of predictions, especially when class distributions are relatively balanced. In contrast, F1-score,



which harmonizes precision and recall, is more suitable for scenarios with class imbalance, a common challenge in log analysis where normal logs significantly outnumber anomalies. Although metrics like ROC-AUC (Fawcett 2006) and PR-AUC were considered, they were ultimately excluded due to their reliance on probabilistic outputs and threshold variation, which are not directly applicable to classification-style outputs generated by LLMs through prompt engineering. For tasks such as log parsing and log summary, where multiple valid outputs may exist, we use ROUGE-L (Lin 2004) and Edit Distance to measure semantic and structural similarity between the generated and reference texts. ROUGE-L evaluates the longest common subsequence between two texts, capturing the structural overlap, which is ideal for assessing key information extraction in templates or summaries. Edit Distance quantifies the number of character-level operations needed to transform the generated output into the reference text, making it especially useful for tasks with strict format constraints such as log parsing. We also considered BLEU (Papineni et al. 2002), a common n-gram based metric, but it was not chosen due to its sensitivity to word order and reduced robustness in tasks with high output variability like summarization and parsing. To comprehensively assess the efficiency and usability of LLMs, we additionally incorporate Average Number of Tokens and Inference Time. These two metrics are critical for realworld applications-longer outputs often imply higher computational and memory costs, and longer inference times directly affect system responsiveness. While alternative systemlevel metrics were considered, they were excluded due to cross-platform inconsistency and difficulty in reproducibility. Instead, token count and time are universally measurable and meaningful across different LLMs and environments. The metrics and formulas presented above help evaluate the performance of each task, ensuring that the system can quickly identify issues, diagnose faults, and generate actionable summaries for system operators.

#### 3 Related Work

In this section, we first discuss the existing evaluations of LLMs in general NLP tasks, as our research also focuses on evaluating LLMs. These evaluations highlight the broad applications of LLMs in NLP. However, there is currently no systematic evaluation of LLMs specifically in the field of log analysis. Therefore, we also examine the applications of LLMs in log analysis tasks, providing context for our research and emphasizing the potential of LLMs in this area, as well as the current lack of standardized evaluation frameworks.

#### 3.1 Evaluation of LLMs in General NLP Tasks

The evaluation of LLMs in general NLP tasks has diversified, as these models have become capable of handling increasingly complex and varied tasks. Such evaluations now not only measure basic linguistic understanding and generation, but also delve into nuanced capabilities such as reasoning, adaptability to different tasks, and the use of domain-specific knowledge. We categorize these evaluations into two main areas: general domain evaluations and specialized domain evaluations.

**Evaluations in General Domain** Comprehensive assessments are designed to evaluate the broad capabilities of LLMs across multiple dimensions. For instance, HELM (Liang et al.



2022) utilizes a diverse set of metrics to assess LLMs in 42 unique scenarios, providing insights into their general linguistic abilities and reasoning skills. BIG-bench (Srivastava et al. 2022) extends this by including tasks that challenge the models' understanding of common sense, logic, and even creativity.

Evaluations In Specialized Domains These assessments focus on evaluating LLMs' performance in domains requiring specialized knowledge. For example, FinEval (Zhang et al. 2023) measures financial acumen, while MultiMedQA (Singhal et al. 2023) tests medical knowledge using datasets derived from professional exams and consultation records. Similarly, Huatuo-26M (Li et al. 2023) evaluates medical consultation capabilities, reflecting real-world medical inquiry handling. NetOps (Miao et al. 2023) focuses on network operations, and tests LLMs with tasks that mimic real-world challenges in network management. OpsEVAL (Liu et al. 2023) assesses the ability of LLMs to manage IT operations, through a set of structured tasks, in both Chinese and English. RepairBench (Silva and Monperrus 2024) establishes an execution-based leaderboard for program repair, evaluating LLMs on real-world Java bugs through test-suite validation and syntactic analysis, providing standardized assessment for AI-driven code repair.

## 3.2 Applications of LLMs in Log Analysis Tasks

With the growing application of LLMs in log analysis tasks, researchers are increasingly exploring how these models can improve key processes such as log parsing and anomaly detection.

Log Parsing LILAC (Jiang et al. 2024) leverages the in-context learning (ICL) capabilities of LLMs by employing a hierarchical candidate sampling algorithm to select high-quality examples for log template generation. It also introduces an adaptive parsing cache to store and refine templates generated by LLMs, reducing query frequency and ensuring template consistency. LogParser-LLM (Zhong et al. 2024) combines the semantic understanding capabilities of LLMs with a prefix tree clustering approach. It utilizes LLMs to process the semantic information of logs and performs online log parsing without requiring hyperparameter tuning or labeled data. DivLog (Xu et al. 2023) uses the ICL capabilities of LLMs to select diverse offline log samples as candidate examples, it then constructs prompts to generate log templates for target logs, enabling unsupervised log parsing. ECLIPSE (Zhang et al. 2024) integrates the semantic understanding capabilities of LLMs with data-driven template matching algorithms to handle cross-lingual log parsing. LLMs are used to extract semantic information from log keywords, improving parsing efficiency in cross-lingual environments. LogPrompt (Liu et al. 2024) employs the zero-shot capabilities of LLMs and advanced prompt strategies to perform log parsing tasks, it enhances LLM interpretability and flexibility, enabling log analysis without relying on training data.

Log Anomaly Detection LogExpert (Wang et al. 2024) integrates LLMs with domain knowledge from technical forums such as Stack Overflow. LLMs are utilized to parse relevant technical solutions and automatically generate recommended resolutions for anomalous logs, reducing the need for manual intervention. SeaLog (Liu et al. 2023) employs LLMs, such as ChatGPT, to provide expert-level feedback that enhances the accuracy of its



Trie-based Detection Agent (TDA) for real-time anomaly detection, allowing the system to adapt to evolving log data more effectively. LogGPT (Qi et al. 2023) utilizes ChatGPT's language understanding and knowledge transfer capabilities through prompt-based techniques for log anomaly detection, exploring the application of large-scale corpora knowledge to the processing of complex log data. LogPrompt (Liu et al. 2024) everages the zero-shot capabilities of LLMs through a set of advanced prompting strategies specifically designed for log anomaly detection tasks. This approach enables LLMs to perform detection without relying on training data, while also offering interpretability of the results.

Other Applications In addition to log parsing and anomaly detection, LLMs have potential applications in various aspects of log analysis. For example, Face It Yourselves (Shan et al. 2024) introduces a two-stage, LLM-based framework for diagnosing configuration errors through log analysis. This framework, called LogConfigLocalizer, leverages LLMs to help endusers, particularly those without source code access, identify the root causes of configuration issues by analyzing logs. UniLog (Xu et al. 2024) employs the ICL paradigm of LLMs to automatically generate appropriate log statements without requiring any fine-tuning. By using prompts with a few demonstration examples, LLMs can determine log positions, verbosity levels, and generate log messages, thus aiding in system maintenance and troubleshooting. LLM4Sec (Karlsen et al. 2024) utilizes various LLM architectures, such as BERT, RoBERTa, and GPT-2 (Radford et al. 2019), to analyze log files for cybersecurity purposes. These LLMs are fine-tuned for specific log types to enhance security log analysis. Summary Cycles (Block et al. 2023) applies LLMs, specifically ChatGPT, to summarize interaction logs in collaborative intelligence analysis. LLMs are used iteratively with recursive summarization techniques to extract key entities, topics, and summaries from user interaction sequences.

However, currently there is no dedicated benchmark for evaluating the performance of different LLMs in various log analysis tasks. This work bridges this gap and proposes an evaluation framework for LLMs in log analysis. Our evaluation efforts intend to understand the strengths and limitations of different LLMs in various log analysis tasks, while providing valuable resources and guidance for the log analysis domain, promoting the effective application of LLMs in real-world scenarios.

# 4 LogEval Benchmark

In this section, we first introduce the platform architecture and the technical stack behind *LogEval*, which provide the foundation for its operation and scalability. The architecture is designed to ensure flexibility and extensibility, supporting a wide range of log analysis tasks. Following this, we describe the key components of our benchmark and its specific evaluation process.

# 4.1 Platform Architecture and Technology Stack

The design and implementation of the *LogEval* platform rely on a powerful and flexible technology stack that ensures high scalability, efficient processing, and easy extensibility. Below, we highlight the core aspects of the platform's architecture and the tools chosen for log analysis.



#### 4.1.1 Platform Architecture

The *LogEval* platform is primarily developed using Python 3.9.6 and runs on Amazon EC2 servers. Flask 3.0.3 is used for building API interfaces, enabling the platform to handle concurrent requests efficiently. The modular architecture allows components to be updated or replaced without disrupting the entire system, ensuring the platform's scalability. Key features of the architecture include:

- Multilingual Support: Integration of Flask-Babel 4.0.0 enables bilingual support (Chinese and English).
- Flexible Data Access: The platform uses the json 2.0.9 package for Managing data in JSON format and Pandas 2.2.2 for data processing. These tools are chosen for their robustness and ability to handle large-scale data efficiently.

## 4.1.2 Extensibility and Scalability Design

To ensure flexibility, the *LogEval* platform is designed to scale and integrate with new features. The following mechanisms support its extensibility:

- Plugin Mechanism: Users can easily integrate new LLM models or log processing techniques by adding custom plugins. This allows for seamless adaptation to future requirements.
- Modular Architecture: The platform's core functionalities, such as log parsing and
  fault diagnosis, are designed as independent modules. New modules can be added as
  needed without modifying the underlying system.
- API Interfaces: The platform provides open API interfaces to enable users to integrate
  with external systems and extend functionality. For example, new LLM models can be
  integrated via simple API calls, allowing users to switch models based on task requirements.
- Hardware Configurations for Performance Testing: The platform's performance across various tasks May be influenced by the underlying hardware configurations. For local deployments, the platform uses a high-performance setup, including 8 NVIDIA A6000 GPUs, each equipped with 48GB of memory, and Intel Xeon processors. For external API calls, the platform uses the official recommended API interfaces provided by the API provider, ensuring consistency and fairness in performance evaluations.

This scalable and modular design ensures that *LogEval* can adapt to future needs, whether it involves adding new features, models, or data sources.

#### 4.2 Evaluation Benchmark

In this section, we introduce the evaluation benchmark *LogEval*, which is designed to assess the performance of various LLMs in performing log analysis tasks. As shown in Fig. 2.



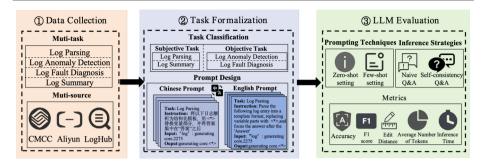


Fig. 2 The framework of LogEval

#### 4.2.1 Data Collection

Data Collection is the foundational step that supports the entire evaluation process. To ensure comprehensive assessment, we curated datasets from diverse sources and tasks, covering a wide range of log processing needs. We designed four core tasks to evaluate LLM capabilities across different log analysis scenarios. In addition, we integrated multi-source datasets to enhance the framework's adaptability and generalizability:

- Aliyun: The dataset contains a total of 299,817 logs, which are grouped by serial number and sorted chronologically. The dataset captures three main fault categories, and their root causes, as flagged by operation and maintenance staff, include issues such as high CPU temperature, memory leaks, and hardware crashes. The dataset provides a real-world perspective on server failures, enhancing its value for anomaly detection and fault diagnosis tasks. The dataset is publicly available at <a href="https://tianchi.aliyun.com/competition/entrance/531947/information">https://tianchi.aliyun.com/competition/entrance/531947/information</a>.
- CMCC: The dataset consists of 482,515 logs collected from OpenStack's (Rosado and Bernardino 2014) OpenVSwitch services, distributed across 493 nodes in a high-performance computing cluster. This dataset spans six fault categories, with root causes ranging from software bugs to resource underprovisioning and unexpected process restarts. The dataset's scale and complexity, derived from an industrial OpenStack environment, make it an excellent benchmark for evaluating anomaly detection methods. The dataset is available at https://github.com/SycIsDD/LogKG.
- LogHub: (Jiang et al. 2024) This open-source repository contains large-scale logs from
  multiple open-source projects, covering real-world scenarios in industries such as server
  management and cloud computing. These datasets not only feature extensive diversity
  but also include detailed annotations, providing a reliable foundation for evaluating the
  performance of LLMs in log processing tasks.

By combining multi-task and multi-source datasets, the *LogEval* framework simulates real-world production environments, providing a solid foundation for comprehensive evaluation of LLM performance in log processing.



#### 4.2.2 Task Formalization

The purpose of this step is to structure the log analysis tasks to match the input requirements of LLMs, thereby achieving effective LLM evaluation and comparison. Task classification is a core step in the formalization process. Based on the nature of task evaluation, we categorize log analysis tasks into two main types:

- Subjective Tasks: These include Log Parsing and Log Summary. These tasks do not
  have a unique correct answer and rely on semantic understanding and content generation for assessment.
- Objective Tasks: These include Log Anomaly Detection and Log Fault Diagnosis.
   These tasks have definite answers, allowing for straightforward quantitative evaluation.

Prompt design is another key aspect to ensure that LLMs can understand and effectively complete each task, each prompt consists of the following four elements:

- Task: Clearly specifies the log analysis task to be evaluated, such as Parsing (Log Parsing), Detection (Log Anomaly Detection), Diagnosis (Log Fault Diagnosis), and Summary (Log Summary).
- **Instruction:** Thoroughly describes the task requirements, guiding the LLM's behavior, for example, instructing the LLM on how to transform a log entry into a structured format.
- **Input:** Provides the log entry or sequence to be analyzed, presented in a uniform format, prefixed with explicit labels like "log entry:" or "log sequence:".
- **Output:** Defines the format of the response to ensure that the LLM's output meets the expected standards.

To evaluate LLMs' performance across different languages, we have prepared prompts in both Chinese and English for each task. Additionally, we provide each task with 15 different prompts to minimize the influence of prompt variations. Table 1 gives three different English prompts for each task.

Table 1 Three Different English Prompts for Each Task

Tasks	English Prompts
Log Parsing	1. Parse the following log into a template format, replacing variable parts with <*>: [log] 2. Convert the following log into a standardized template by identifying and replacing the variable parts with <*>: [log] 3. Transform the raw log [log] into a log template by replacing variable segments with <*>
Log Anomaly Detection	Review and mark the log entry as "normal" or "anomalous", only output "normal" or "anomalous"     Analyze the log content, classify it as "normal" or "anomalous", only output "normal" or "anomalous"     Check the log entry, and determine if it belongs to the "normal" or "anomalous" category, only output "normal" or "anomalous"
Log Fault Diagnosis	1. In our data scenario, there are several types of faults {fault types}. Analyze the log [log] and identify the type of fault that occurred. Only output the fault type 2. In our data scenario, there are several types of faults {fault types}. Based on the information in the log [log], determine which type of fault the log represents. Only output the fault type 3. In our data scenario, there are several types of faults {fault types}. Use the detailed information provided by the log [log] to conduct an in-depth analysis to determine the category of the fault. Only output the fault type
Log Summary	1. Analyze the following 20 logs [log], extract key information, phrases, sentences, or recurring content to generate a summary, and only output the summary 2. Extract the most important events, phrases, and activities or recurring content from the following 20 logs [log], create a concise log overview, only output the summary 3. Extract key events, sentence phrases, or recurring information from the following 20 logs [log] to form a comprehensive summary, only output the summary



#### 4.2.3 LLM Evaluation

labelsubsec:evaluation This section evaluates LLMs' capabilities in log analysis through systematic benchmarking. We first introduce the evaluation strategies, then detail the selected models. Our benchmarking framework combines two evaluation strategy:

**Inference Strategy** We employ two different inference strategies to process and interpret the responses generated by LLMs: Naive Q&A and Self-Consistency Q&A. These strategies aim to investigate the stability of LLM outputs.

- Naive Q&A: This strategy involves a single model invocation per query, and the generated answer is directly treated as the final prediction. Naive Q&A is simple and efficient, and it is especially suitable for tasks with subjective nature and diverse valid answers, such as log parsing and summarization.
- Self-Consistency Q&A: To enhance the stability and accuracy of model outputs, Self-Consistency Q&A performs multiple model invocations on the same query (set to 5 times in our study), generating multiple answers. The most frequent answer among these is selected as the final result through a voting mechanism. This approach effectively reduces the randomness of single-shot outputs and is particularly well-suited for tasks with objective ground truth, such as log anomaly detection and fault diagnosis.

**Prompting Technique** We use various settings to evaluate LLMs on LogEval to get a comprehensive overview of their performance. We evaluate LLMs in zero-shot and few-shot (5-shot) settings.

- **Zero-shot setting**: This technique involves presenting the LLM with a task without prior examples, thereby testing its ability to adapt to new situations based on its pre-existing knowledge. It is a measure of the LLM's capacity to generalize from its training data to unseen tasks. The examples for the zero-shot setting can be found in Table 2.
- *Few-shot setting*: The LLM is provided with a limited number of exemplars before being asked to perform the task. Few-shot prompting helps the model better capture task-specific patterns or structures within the log data, often leading to improved performance compared to zero-shot. The examples for the few-shot setting can be found in Table 3.

Table 2 Zero-sho	prompts f	for the four	log analysis tasks

Task	Parsing	Anomaly Detcetion	Diagnosis	Summary
Instruction	Parse the following log entry into a template format, replacing variable parts with <*>, and focus the answer after the keyword 'Answer'	Please review the log entry and explicitly mark it as 'normal' or 'anomalous', only output 'normal' or 'anomalous'	In our data scenario, there are three types of faults: Processor CPU Cates. Memory Throttled Uncorrectable Error Correcting Code, Hard Disk Drive Control Error Computer System Bus Short Circuit Programmable Gate Array Device Unknown. Analyze the log entry and identify the type of fault that occurred. Only output the fault type.	Analyze the following 20 logs, extract key information, phrases, sentences, or recurring content to generate a summary, only output the summary.
Input	log entry: synchronized to 10.100.28.250, stratum 3	log entry: instruction cache parity error corrected	log entry: Processor #0xfa   Configuration Error Asserted	log sequence:[blockMap updated: 10.251.193.175:50010 is added to blk 3864576029521084501 size 3540711,
Output	synchronized to <*>, stratum <*>"	normal	Processor CPU Caterr	blockMap updated; PacketResponder terminating; Received block;



**Table 3** Few-shot prompts for the four log analysis tasks

Task	Parsing	Diagnosis
Instruction	Parse the following log entry into a template format, replacing variable parts with <>>, and focus the answer after the keyword 'Answer' For example:log entry; no floppy controllers found, answer: "no floppy controllers found, in gentry: 13 tree receiver in re-synch state event(6) deteroid ver 4-55 decends, answer: '<>> in re-synch state event(6) (der <>>) detected over <>> seconds'; log entry: a sutorna DONE; log entry: a sutorna DONE; log entry: a sutorna DONE; log entry: a L3 EDRAM error(5) (der 'do, 157) detected and corrected over 282 seconds, snew: '<>> L5 EDRAM error(6) (der <>> detected and corrected over <pre>cected ov</pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre>	In our data scenario, there are three types of faults (Processor CPU Caterr, Memory Throttled Uncorrectable Error Correcting Code, Hard Disk Drive Control Error Computer System Bus Short Circuit Programmable Gate Array Device Unknown). Analyze the log entry and identify the type of fault that occurred. Only output the fault type. For Example: log entry: Temperature CPU0. Margin, Temp   Lower Critical going low   Asserted   Reading : 16 = Threshold 0 degrees C answer: Processor CPU Caterr'; log entry: Memory CPU1ED. DIMM. Stat   Correctable ECC   Asserted answer: "Memory Throttled Uncorrectable Error Correcting Code'; log entry: System Boot Initiated BIOS. Boot. Up   Initiated by power up   Assertedanswer: "Hard Disk Drive Control Error Computer System Bus ShortCircuit Programmable Gate Array Device Unknown."
Input	log entry: synchronized to 10.100.28.250, stratum 3	log entry: Processor #0xfa   Configuration Error  Asserted
Output	synchronized to <*>, stratum <*>*	Processor CPU Caterr

Table 4 LLMs Chosen for Evaluation

Model	Creator	Size	Access
GPT-4 OpenAI et al. 2024	OpenAI	undisclosed	Commercial
GPT-3.5 OpenAI 2022	OpenAI	undisclosed	Commercial
Claude-3-Sonnet Anthropic 2023	Anthropic	undisclosed	Commercial
Gemini-Pro Team et al. 2023	Google	undisclosed	Commercial
Mistral Jiang et al. 2023	Mistral	7B	Open-source
InternLM2-Chat InternLM 2023	Shanghai AI Laboratory	7B/20B	Open-source
DevOps-Model-Chat CodeFuse 2023	CodeFuse	7B/14B	Open-source
AquilaChat BAAI 2023	BAAI	7B	Open-source
ChatGLM-4 THUDM 2024	Tsinghua	undisclosed	Commercial
LLaMA-2 Touvron et al. 2023	Meta	7/13/70B	Open-source
Qwen-1.5-Chat Bai et al. 2023	Alibaba Cloud	7/14/72B	Open-source
Baichuan2-Chat Yang et al. 2023	Baichuan Intelligence	13B	Open-source

We select 12 state-of-the-art LLMs covering diverse architectures and accessibility modes, as summarized in Table 4.

#### 5 Evaluation Results

In this section, we aim to explore the following key aspects of LLMs' performance in log analysis tasks:

- **RQ1:** What is the overall performance of different LLMs when applied to various log analysis tasks?
- RQ2: How do LLMs perform under Naive Q&A settings across different log analysis tasks?
- **RQ3:** How do LLMs perform under Self-Consistency Q&A settings in the context of log analysis tasks?
- **RQ4:** What is the impact of inference time and the average number of tokens on the performance of LLMs?
- RQ5: How do factors such as parameter size and language choice influence the performance of LLMs in log analysis tasks?



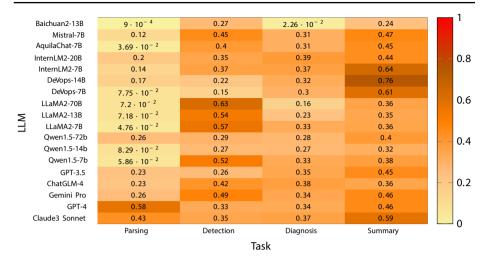


Fig. 3 Accuracy in zero-shot Naive Q&A across four tasks

						_	. 1
	Baichuan2-13B	2.5 · 10 <sup>-4</sup>	0.22	1.5 · 10 <sup>-2</sup>	0.42		l '
	Mistral-7B	0.14	0.63	0.68	0.8		
	AquilaChat-7B	3.59 · 10 <sup>-2</sup>	0.25	0.13	0.19		
	InternLM2-20B	0.59	0.34	0.61	0.64		0.8
	InternLM2-7B	0.31	0.32	0.58	0.67		
	DevOps-14B	0.17	0.26	0.52	0.85		
	DevOps-7B	0.1	0.2	0.47	0.81		
	LLaMA2-70B	8.3 · 10 <sup>-2</sup>	2.09 · 10 <sup>-2</sup>	0.3	0.44		0.6
>	LLaMA2-13B	2.23 · 10 <sup>-2</sup>	2.5 · 10 <sup>-4</sup>	6.08 · 10 <sup>-2</sup>	0.19		
LLM	LLaMA2-7B	8.27 · 10 <sup>-2</sup>	2.13 · 10 <sup>-3</sup>	4.89 · 10 <sup>-2</sup>	0.21		
	Qwen1.5-72B	0.62	0.3	0.83	0.78		0.4
	Qwen1.5-14B	0.43	7.05 · 10 <sup>-2</sup>	0.5	0.68		0. 1
	Qwen1.5-7B	0.27	4.95 · 10 <sup>-2</sup>	0.52	0.29		
	GPT-3.5	0.23	0.39	0.67	0.47		
	ChatGLM-4	0.57	0.36	0.71	0.56		0.2
	Gemini Pro	0.84	0.57	0.83	0.73		
	GPT-4	0.89	0.53	0.91	0.54		
	Claude3 Sonnet	0.87	0.31	0.67	0.53		0
		Parsing	Detection	Diagnosis	Summary		. 0
			Ta	sk			

Fig. 4 Accuracy in few-shot Naive Q&A across four tasks

#### 5.1 RQ1: Overall Performance

To evaluate the performance of various LLMs on different log analysis tasks, we conducted a comparative analysis of their Naive Q&A accuracy under both zero-shot and few-shot settings. The results are shown in Figs. 3 and 4, respectively. For the sake of simplicity, we use the abbreviation of each task in these two and subsequent figures, *i.e.*, we use "Parsing" instead of "Log Parsing", "Detection" instead of "Log Anomaly Detection", "Diagnosis" instead of "Log Fault Diagnosis", and "Summary" instead of "Log Summary". From Figs. 3 and 4, we have the following findings for each task:



- Log Parsing: For log parsing, GPT-4 and Claude3 Sonnet demonstrate outstanding performance in both zero-shot and few-shot settings, with GPT-4 achieving the highest parsing accuracy under the few-shot condition, showcasing its exceptional parsing capabilities.
  Gemini Pro also exhibits strong adaptability in the few-shot setting, achieving a high level
  of parsing accuracy, which positions it as a competitive and promising LLM for this task.
- Log Anomaly Detection: In the log anomaly detection task, LLaMA2-70B performs better than other LLMs in the zero-shot setting, but it still lags slightly behind GPT-4 and Claude3 Sonnet in overall performance. In the few-shot setting, Mistral-7B shows a significant improvement, demonstrating strong contextual learning abilities, making it the standout LLM in this task. Gemini Pro also performs well in the few-shot setting, showcasing its adaptability to different prompt conditions, making it suitable for applications in dynamic data environments.
- Log Fault Diagnosis: In the log fault diagnosis task, performance differences among LLMs in the zero-shot setting are relatively small; however, Baichuan2-13B and the LLaMA2 series LLMs show relatively weaker performance in this task. In the few-shot setting, GPT-4 shows a marked improvement, establishing itself as the best choice for this task, while Gemini Pro and Qwen1.5-72B also exhibit excellent diagnostic capabilities. These results suggest that GPT-4 can effectively enhance fault diagnosis accuracy under few-shot conditions, making it an ideal LLM for complex diagnostic tasks.
- Log Summary: In the log summarization task, the DeVops series LLMs perform well in both zero-shot and few-shot settings, showing their advantage in summary generation.
   In the few-shot setting, Mistral-7B and Qwen1.5-72B show significant improvement, demonstrating the ability to generate high-quality log summaries with limited prompts.
   These LLMs have application potential in scenarios requiring high-quality log summarization, especially where limited data is available.

We further compare the accuracy of Commercial and Open-source LLMs, in accordance with the access types listed in the "Access" column of Table 4. The results are shown in Figs. 5 and 6, from which we can draw the following key findings:

- Zero-shot Setting Analysis: In zero-shot scenarios, weight-based LLMs generally outperform API-based LLMs, with InternLM2-20B and Mistral-7B standing out for their high accuracy, demonstrating the stability and superior performance of weight-based LLMs in local runtime environments. Among the API-based LLMs, Claude3 Sonnet and GPT-4 show relatively stable performance, indicating that in multi-task scenarios, these LLMs can deliver reliable performance under zero-shot conditions, making them suitable for generic task applications that do not require fine-tuning.
- Few-shot Setting Analysis: In few-shot settings, weight-based LLMs show significant performance improvements, with InternLM2-20B and Mistral-7B exhibiting high adaptability with few-shot prompts. API-based LLMs also see noticeable improvement in the few-shot setting, especially with Gemini Pro and GPT-4 achieving high few-shot accuracy, demonstrating strong adaptability. However, weight-based LLMs demonstrate a more pronounced capacity for adaptation in few-shot learning, making them well-suited for complex task scenarios requiring frequent updates and optimizations. In contrast, API-based LLMs, with limited fine-tuning flexibility, are better suited for applications requiring stability and immediate responsiveness.



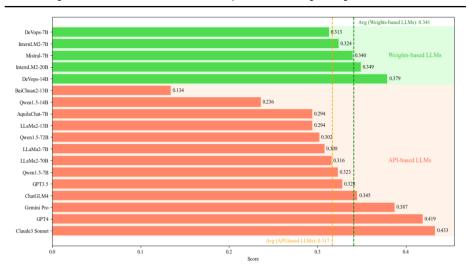


Fig. 5 Overall Performance in zero-shot Naive Q&A

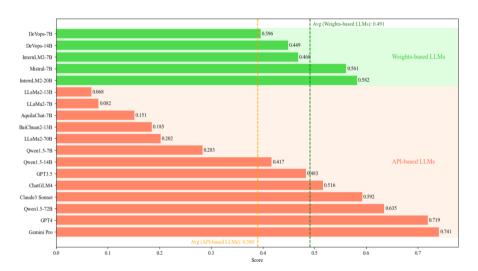


Fig. 6 Overall Performance in few-shot Naive Q&A

#### 5.2 RQ2: Naive Q&A Performance

To investigate the performance of various LLMs in Naive Q&A across different log analysis tasks, we conducted a comparative analysis of their performance under both zero-shot and few-shot settings. This section examines the results for each task, highlights the strengths and weaknesses of different models, and explores the potential factors influencing their performance,we present the following findings for each task.



## 5.2.1 Naive Q&A results on Log Parsing

We evaluated the performance of various LLMs on Naive Q&A log parsing task in both zero-shot and few-shot settings. The following conclusions can be drawn:

- Few-shot learning consistently boosts LLM accuracy: Across most LLMs, few-shot
  learning substantially improves accuracy compared to zero-shot settings. This improvement is particularly notable in high-performing LLMs such as GPT-4 and Claude3Sonnet, indicating that few-shot learning can effectively enhance LLM adaptability to
  complex log parsing tasks.
- GPT-4 and Claude3-Sonnet excel in multiple parsing tasks: Among the evaluated LLMs, GPT-4 and Claude3-Sonnet consistently deliver high performance across both Chinese and English log parsing tasks in zero-shot and few-shot settings. Their robust accuracy and low Edit Distance suggest strong generalization and adaptability across languages and parsing scenarios.
- LLM performance scales with LLM size and architecture: The performance data reveals that larger, more sophisticated LLMs, such as GPT-4 and Claude 3-Sonnet, consistently outperform smaller LLMs, including BaiChuan2-13B and AquilaChat-7B. This scaling effect underscores the advantage of larger LLMs with advanced architectures in capturing complex patterns in log parsing tasks, while smaller LLMs struggle to generalize and adapt effectively.

## 5.2.2 Naive Q&A results on Log Anomaly Detection

We evaluated the performance of various LLMs on Naive Q&A log anomaly detection task in both zero-shot and few-shot settings. The following conclusions can be drawn:

- Limited Impact of Few-shot Learning: In both Chinese and English tasks, few-shot learning does not significantly outperform zero-shot learning in log anomaly detection. This may be due to the LLMs' tendency to reproduce the answers found in the few-shot examples, thereby failing to effectively identify new anomaly cases. Table 5 illustrates some specific examples of anomaly detection tasks where the LLM repeats answers given in the samples, reflecting the LLM's performance in real-world applications. This indicates that simply increasing the number of samples may not significantly enhance LLM performance, especially when dealing with imbalanced datasets where anomaly classes are underrepresented.
- Sensitivity to Prompt Language: The differences in performance between Chinese and English prompts are not substantial overall, but some LLMs do exhibit varying performance levels depending on the prompt language. For example, certain LLMs, such as LLaMA2-70B, show lower F1-Score under Chinese prompts, particularly in the few-shot setting, where an F1-Score of 0 was observed. This suggests that these LLMs have limited generalization capabilities when processing Chinese prompts, highlighting the sensitivity of the LLM to the prompt language.



Table 5 Few error examples on Log Anomaly Detection

"Task": Detection	"Task": Detection	"Task": Detection	"Task": Detection
"label": "normal"	"label": "normal"	"label": "normal"	"label": "normal"
"predict": "log1: normal log2: normal log3: anomalous 4: anomalous 5: normal 6 anomalous"	"predict": "normal normal anomalous normal"	"predict": "'normal' 'normal' 'abnormal'"	"predict": "normal normal abnorma anomalous"

Table 6 A failed zero-shot example by BaiChuan2-13B

"Task": "Detection"	"Task": "Detection"		
"label": "anomalous"	"label":"normal"		
"predict": "anomalous"	"predict": "Based solely on the information provided in the log entry, it is not possible to definitively classify it as 'normal' or 'anomalous'. The log entry mentions "double-hummer alignment exceptions", but without additional context or information about what this means within the specific system or application being monitored, it is difficult to determine whether this is expected behavior or an indication of a problem. Therefore, I would classify this log entry as 'unclear' or 'inconclusive' rather than 'normal' or 'anomalous'. It may be necessary to gather more information or consult additional logs or documentation in order to make a more accurate classification."		

• Challenges with Complex or Domain-Specific Questions: In zero-shot settings, LLMs often struggle with addressing complex or domain-specific questions, resulting in vague or uncertain outputs. As illustrated by the BaiChuan2-13B model's performance on log analysis tasks (Table 6), even high-performing models may fail to accurately classify log entries without sufficient domain knowledge. Integrating domain-specific information into the training process can significantly improve comprehension and response accuracy for specialized tasks like log analysis.

#### 5.2.3 Naive Q&A Results on Log Fault Diagnosis

We evaluated the performance of various LLMs on Naive Q&A log fault diagnosis task in both zero-shot and few-shot settings. The following conclusions can be drawn:

Effectiveness of Few-shot Learning: Few-shot learning markedly enhances accuracy
and F1-Score across most LLMs. High-performing LLMs, such as GPT-4 and Qwen1.572B, show significant improvements in the few-shot setting, highlighting the value of
providing examples in fault diagnosis. However, some smaller LLMs, like the LLaMA
series, exhibit limited benefits from few-shot learning, indicating their difficulty in



- adapting to complex tasks through minimal examples.
- Superiority of GPT in Complex Tasks: Among the evaluated LLMs, GPT-3.5 and GPT-4 perform exceptionally well in few-shot log fault diagnosis, with both LLMs achieving an F1-Score above 0.9. GPT-4's superior performance in both Chinese and English tasks suggests that it can effectively generalize in complex fault diagnosis scenarios, making it ideal for critical applications requiring high accuracy.
- Limitations of Smaller LLMs: Smaller LLMs with fewer parameters exhibit significant limitations in handling complex fault diagnosis tasks. Their reduced capacity to capture diverse fault types, particularly in zero-shot settings, leads to inaccurate predictions. Table 7 showcases several failure examples from 7B models. For instance, when tasked with diagnosing "Processor CPU Caterr," the model inaccurately predicts "Processor CPUR," and for "Mirror does not specify output," the model erroneously outputs "Mirror specified output." These issues demonstrate that smaller models struggle to capture the deeper, intricate features of log data, resulting in lower reliability of their fault diagnosis predictions.

## 5.2.4 Naive Q&A Results on Log Summary

We evaluated the performance of various LLMs on Naive Q&A log summary task in both zero-shot and few-shot settings. The following conclusions can be drawn:

- Effectiveness of Few-shot Learning: The majority of LLMs show significant improvement in log summarization performance when transitioning from zero-shot to few-shot settings. This suggests that few-shot learning enables LLMs to better grasp and adapt to the summarization context, enhancing both accuracy and F1-Score. However, certain LLMs, especially in the LLaMA series and AquilaChat-7B, display a counterintuitive trend, with performance sometimes higher in zero-shot than in few-shot settings, possibly due to noise introduced by few-shot examples.
- Strong Performance of DeVops: The DeVops-7B and DeVops-14B models consistently outperform others across both zero-shot and few-shot settings. DeVops-14B, in particular, demonstrates exceptional summarization capabilities, making it well-suited for applications where accuracy and robustness in log summarization are critical.

Table 7 Few error examples on Log Fault Diagnosis

"Task": Diagnosis	"Task": Diagnosis	"Task": Diagnosis
"label": "Processor CPU Caterr"	"label": "Mirror does not specify output"	"label": "Processor CPU Caterr"
"predict": "answer:'Processor CPUR';"	"predict": "Mirror specified output"	"predict": "1. Processor CPU Cater 2. Memory Throttled Uncorrectable Error Correcting Code 3. None 4. Hard Disk Drive Control Error Computer System Bus Short Circuit Programmable Gate Array Device Unknown"



Task-Specific LLM Performance Diversity: In different tasks and language settings, specific LLMs exhibit notable performance variations, highlighting their adaptability and Limitations in particular tasks or languages. For example, Gemini Pro performs exceptionally well in few-shot English tasks, demonstrating high adaptability, but shows weaker performance in zero-shot Chinese tasks. Similar trends are observed in LLMs Like Claude 3-Sonnet. These results suggest that variations in LLM performance across tasks may reflect the impact of optimization focus and training data.

## 5.3 RQ3: Self-Consistency Q&A Performance

To evaluate the capability of various LLMs in Self-Consistency Q&A for log anomaly detection and log fault diagnosis, as well as self-consistency in LLM robustness performance, we conducted experiments under zero-shot and few-shot settings, and provide a corresponding analysis of these findings, we present the following findings for each task.

## 5.3.1 Self-Consistency Q&A results on Log Anomaly Detection

From the overall performance results, we can draw the following scientifically conclusions:

- Few-shot learning does not outperform zero-shot learning in log anomaly detection tasks, highlighting its limitations in this context. In the Self-Consistency Q&A test, which involves multiple inquiries to the LLM and taking the most frequent answer, few-shot learning did not significantly surpass zero-shot learning. This outcome may be because the provided few-shot examples still do not sufficiently cover all patterns, thus failing to improve LLM consistency. LLMs in this setup tend to repeat examples rather than effectively learn new anomaly detection patterns from limited samples.
- The BaiChuan model shows a significant improvement in the Self-Consistency mode, indicating potential for more consistent responses, though its performance remains volatile. Compared to the Naive Q&A test, the BaiChuan model improved notably in the Self-Consistency Q&A test, suggesting a greater likelihood of generating consistent answers in repeated queries. However, it also shows considerable variability in responses across rounds, revealing a lack of stability in multi-turn interactions. Further optimization is needed to enhance the BaiChuan model's consistency and coherence in continuous query settings.
- The LLaMA2 series of models demonstrates poor performance and lack of stability in Self-Consistency Q&A test, suggesting the need for further improvements and optimizations. In multiple queries, the LLaMA2 models continue to produce low and inconsistent performance, indicating deficiencies in generating stable responses. This result may stem from limited generalization capabilities in handling complex tasks or a lack of optimization for log anomaly detection tasks. Enhancing the consistency of the LLa-MA2 models in multi-turn Q&A may require architectural improvements or additional fine-tuning on relevant data to improve robustness in repeated queries.



## 5.3.2 Self-Consistency Q&A Results on Log Fault Diagnosis

From the overall performance results, we find that the few-shot results are better than zero-shot results, similar to the Naive Q&A results. This indicates stable output in the log fault diagnosis task, with GPT-3.5 and GPT-4 showing far superior results. The Baichuan model performs poorly under both Self-Consistency and Naive Q&A, while other LLMs do not change much relative to the Naive Q&A results. The zero-shot and few-shot performance of the LLMs are examined for English and Chinese test sets by comparing the results of the Naive and Self-Consistency Q&A experiment. The following conclusions can be drawn from the results:

- For most LLMs, performance does not change much from Naive Q&A to Self-Consistency Q&A. In the anomaly detection task, the performance under few-shot conditions is inferior to zero-shot. Conversely, in the fault diagnosis task, the performance under few-shot conditions exceeds zero-shot scenarios.
- In these settings, Self-Consistency prompts relatively minor improvements to the LLM.
   In repeated questions, the LLM's answers were consistent.

## 5.4 RQ4: Performance on Inference Time and Average Number of Tokens

To investigate the reasoning efficiency of LLMs and whether they are in generating responses, we summarized the inference time for different LLMs and the average number of tokens output per log. The inference time and Average Number of Tokens used for each task on the English dataset in the zero-shot case of the Naive Q&A are shown below.

#### 5.4.1 Inference Time

Figure 7 presents the inference time of 18 mainstream LLMs across four log analysis tasks, measured under the English dataset and zero-shot Naive Q&A setting. We first analyze the inference performance by task and model, and then discuss their potential in high-throughput scenarios.

Task-wise Inference Time Comparison The log summarization task generally exhibits the longest inference time, with some models reaching 5–7 seconds. This is primarily due to the longer input length and the need for the model to integrate and rewrite information across multiple sentences. Log fault diagnosis and log parsing tasks show moderate inference time (mostly 1–3 seconds), indicating a relatively structured reasoning path and lower computational demand. Log anomaly detection, the only real-time task, achieves the shortest inference time. Lightweight models like DevOps-7B and InternLM2-7B Maintain consistent latency between 0.4–0.7 seconds, demonstrating their suitability for real-time applications.

We also observe a clear correlation between model size and inference latency. 70B-scale models (e.g., LLaMA2-70B, Qwen1.5-72B) show significantly higher latency and are more appropriate for offline tasks, while 7B/14B models provide excellent responsiveness suitable for latency-sensitive deployments.



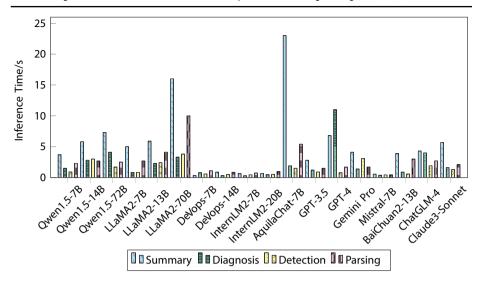


Fig. 7 The Inference Time in the Naive Q&A in log analysis

Analysis of Scalability under High Log Volumes Inference time directly affects a model's capacity to handle large-scale log streams. Based on our measurements, we further analyze the models' applicability in high-throughput industrial scenarios. For example, typical production systems generate approximately 100,000 logs/hour (28 logs/sec). Among the four tasks, only log anomaly detection requires real-time processing. DevOps-7B, with an average latency of 0.43s, can theoretically support over 2,000 logs/sec, exceeding real-world demands and ensuring both low latency and system stability. The remaining three tasks can be processed in offline batches, allowing for the use of larger models (e.g., Qwen1.5-72B) that trade latency for improved accuracy. A practical solution involves a two-stage architecture, Stage 1 (Light Filtering): Rule-based filters or Lightweight LLMs remove 90% of normal logs. Stage 2 (LLM Analysis): The remaining 10% (2.8 logs/sec) are processed by more capable LLMs.

Furthermore, LogGPT (Qi et al. 2023) and LogPrompt (Liu et al. 2024) have demonstrated the ability to process log anomaly detection, further validating the scalability of LLM-based log analysis pipelines. In summary, inference time serves as a practical indicator not only for real-time responsiveness but also for guiding the architectural design of LLM-based solutions to meet high-throughput industrial requirements.

#### 5.4.2 Average Number of Tokens

Figure 8 shows the Average Number of Tokens of the four classes of tasks on the English data set with zero-shot setting for Naive Q&A.

From the overall performance evaluation results, the log summary task outputs the highest average number of tokens among the four tasks. This phenomenon is mainly determined by the nature of the task because the log summary task requires the LLM to generate a concise summary, which usually requires more tokens to accurately represent



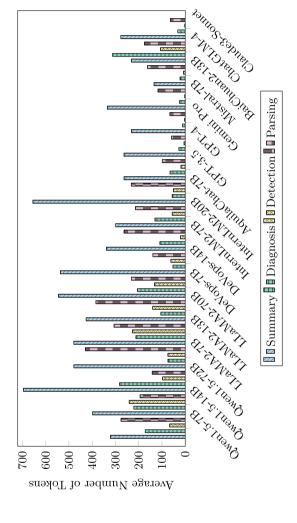


Fig. 8 The Average Number of Tokens in the Naive Q&A in log analysis



the main content of the log. However, our evaluation results show that ChatGLM-4, GPT, and Mistral models output a lower average number of tokens, indicating that their answers are more concise, without excessive redundant information, and their outputs are cleaner. Conversely, LLaMA and Qwen models output more tokens on average, meaning their answers contain more extraneous content. In practice, this can result in users spending more time and effort sifting useful information from responses, which reduces efficiency.

#### 5.5 RQ5: Performance on Different parameters and Language

To provide the impact of different parameter sizes on models, this section conducts a comparative analysis of the performance of LLaMA-2 and Qwen-1.5, each evaluated with three different parameter sizes, offering insights into their adaptability and potential use cases.

Figure 9 shows the accuracy of LLaMA-2 and Qwen-1.5 with different parameter sizes. We used a zero-shot Naive Q&A assessment on English prompts.

From the comparison of results, most LLMs achieve better performance with a parameter size of 7B across the majority of tasks. This finding suggests that LLM size is not a determining factor for log analysis tasks. While an increase in the number of parameters generally means that the LLM can capture more features and patterns, a large number of parameters can also cause the LLM to be too complex to process log data quickly and accurately in real-world applications. Therefore, we can conclude that for log analysis tasks, choosing the right number of parameters is crucial, not simply "bigger is better." Future research should focus on how to optimize the size of the LLM for a more efficient and cost-effective log analysis solution without sacrificing performance.

To provide the impact of different language on models, this section conducts a comparative analysis, as illustrated in Fig. 10, reveals notable differences. LLMs such as LLaMA series, GPT-4, ChatGLM4, and Claude3-Sonnet excel in English tasks, while LLMs like Qwen and DevOps, trained with a substantial amount of Chinese data, outperform in Chinese tasks. This performance disparity is attributable to the linguistic distribution in the LLMs' pretraining datasets. Therefore, task-specific language requirements must guide LLM selection. For Chinese-focused applications, LLMs like Qwen and DevOps are recommended, whereas English-dominant tasks may benefit from the LLaMA series or GPT-4. This discussion outlines the specific performance of LLMs in different languages.

This section provides a comprehensive performance evaluation of several LLMs. Through comparative analysis of these LLMs, we find significant differences in their performance on log analysis tasks. These differences may be due to differences in LLM design philosophy, training strategies, and LLM architecture. For example, some LLMs may perform better in parsing, while others may show greater efficiency in generating summaries or detecting anomalies. Additionally, the number of parameters and training objectives of the LLM are also important factors affecting its performance in the log analysis task. Our evaluation highlights the need to consider these factors when selecting and customizing a log analysis LLM to ensure that the LLM effectively meets the needs of real-world applications.



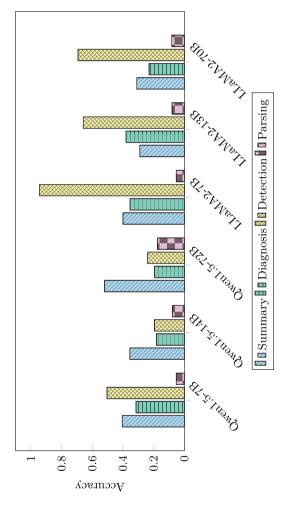


Fig. 9 The Accuracy of LLaMA-2 and Qwen-1.5 in zero-shot English Naive Q&A



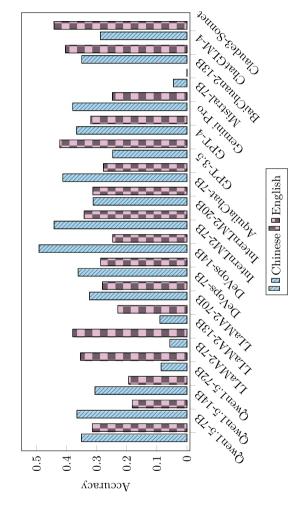


Fig. 10 The performance of LLMs under the "zero-shot" Naive Q&A in both Chinese and English test sets



#### 6 Conclusion

In this study, we have addressed a significant gap in the field of log analysis, where a standardized and systematic evaluation framework for assessing LLMs across multiple tasks has been lacking. The heterogeneity in LLM architectures, varying parameter sizes, and diverse applicability in log analysis complicate the decision-making process for selecting the most suitable LLM. To overcome this challenge, we introduced LogEval, a unified and comprehensive benchmarking suite designed to rigorously evaluate the performance of different LLMs across key log analysis tasks, including log parsing, anomaly detection, fault diagnosis, and summarization. LogEval provides a robust framework that facilitates consistent comparisons among LLMs. The benchmark is complemented by a real-time, dynamically updating platform, accessible at https://nkcs.jops.ai/LogEval/, which serves as a valuable resource for both researchers and practitioners in the domain. This platform enables users to stay up-to-date with the latest advancements in LLM technology and understand how different LLMs perform in practical log analysis scenarios. Our code repository is available at https://github.com/LinDuoming/LogEval. Our evaluation results have highlighted the strengths and limitations of various LLMs, underscoring the importance of task-specific LLM selection and the impact of zero-shot versus few-shot prompting techniques. LogEval not only offers a clear performance overview but also provides insights that can guide the design and deployment of LLM-based log analysis systems.

Moving forward, we plan to expand *LogEval* in three directions: (1) incorporating synthetic and real-world log generation tasks to evaluate LLMs' generative capabilities under structural and semantic constraints; (2) continuously integrating emerging LLMs and instruction-tuning methods to maintain the benchmark's relevance; and (3) collecting more fine-grained industrial logs across diverse domains to support broader downstream evaluations, including security incident response, fault localization, and self-healing automation.

**Acknowledgements** This work is supported by the National Natural Science Foundation of China (62272249, 62302244), and the Fundamental Research Funds for the Central Universities (XXX-63253249).

Author Contributions All authors contributed to the study conception, design, implementation, and manuscript preparation. All authors read and approved the final manuscript.

Funding Not applicable.

Data Availability To facilitate further work by other researchers in this area, all of our scripts and datasets are available online. The material and data from this study are available from the following URL: <a href="https://nkcs.iops.ai/LogEval/">https://nkcs.iops.ai/LogEval/</a>.

## **Declarations**

Ethical approval Not applicable.

Informed consent Not applicable.

Conflict of Interest The authors declare that they have no conflict of interest.

Clinical trial number Not applicable.



#### References

- Cito J, Leitner P, Fritz T, Gall HC (2015) The making of cloud applications: An empirical study on software development for the cloud. In: Proceedings of the 2015 10th joint meeting on foundations of software engineering. ESEC/FSE 2015, Association for Computing Machinery New York NY USA, pp 393–403. https://doi.org/10.1145/2786805.2786826
- Li Y, Jiang ZMJ, Li H, Hassan AE, He C, Huang R, Zeng Z, Wang M, Chen P (2020) Predicting node failures in an ultra-large-scale cloud computing platform: An aiops solution. ACM Trans Softw Eng Methodol 29 (2). https://doi.org/10.1145/3385187
- Zhang X, Xu Y, Qin S, He S, Qiao B, Li Z, Zhang H, Li X, Dang Y, Lin Q, Chintalapati M, Rajmohan S, Zhang D.: Onion: identifying incident-indicating logs for cloud systems. In: Proceedings of the 29th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering. ESEC/FSE 2021, Association for Computing Machinery New York NY USA, pp 1253–1263. https://doi.org/10.1145/3468264.3473919
- Nedelkoski S, Bogatinovski J, Acker A, Cardoso J, Kao O (2020) Self-attentive classification-based anomaly detection in unstructured logs. In: 2020 IEEE international conference on data mining (ICDM), pp 1196–1201. https://doi.org/10.1109/ICDM50108.2020.00148
- Wang J, Chu G, Wang J, Sun H, Qi Q, Wang Y, Qi J, Liao J (2024) Logexpert: Log-based recommended resolutions generation using large language model, pp 42–46. https://doi.org/10.1145/3639476.3639773
- Zhong A, Mo D, Liu G, Liu J, Lu Q, Zhou Q, Wu J, Li Q, Wen Q (2024) Logparser-Ilm: Advancing efficient log parsing with large language models. In: Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining. KDD '24, Association for Computing Machinery, New York NY USA, pp 4559–4570. https://doi.org/10.1145/3637528.3671810
- Locke S, Li H, Chen THP, Shang W, Liu W (2022) Logassist: Assisting log analysis through log summarization. IEEE Trans Softw Eng 48(9):3227–3241. https://doi.org/10.1109/TSE.2021.3083715
- Ma L, Yang W, Xu B, Jiang S, Fei B, Liang J, Zhou M, Xiao Y (2024) Knowlog: Knowledge enhanced pretrained language model for log understanding. In: ICSE, pp 32–13213. https://doi.org/10.1145/35975 03.3623304
- Lin Q, Zhang H, Lou JG, Zhang Y, Chen X (2016) Log clustering based problem identification for online service systems. In: 2016 IEEE/ACM 38th international conference on software engineering companion (ICSE-C), pp 102–111
- Fawcett T (2006) An introduction to roc analysis. Patt Recognit Lett 27(8):861-874
- Meng W, Liu Y, Zaiter F, Zhang S, Chen Y, Zhang Y, Zhu Y, Wang E, Zhang R, Tao S, Yang D, Zhou R, Pei D (2020) Logparse: Making log parsing adaptive through word classification. In: 2020 29th international conference on computer communications and networks (ICCCN), pp 1–9. https://doi.org/10.1109/ICCCN49398.2020.9209681
- He S, He P, Chen Z, Yang T, Su Y, Lyu MR (2021) A survey on automated log analysis for reliability engineering. ACM Comput Surv 54(6):1–37
- Liu Y, Zhang X, He S, Zhang H, Li L, Kang Y, Xu Y, Ma M, Lin, Q, Dang Y, Rajmohan S, Zhang D.: Uniparser: A unified log parser for heterogeneous log data. In: Proceedings of the ACM Web Conference 2022. WWW '22, pp 1893–1901. Association for Computing Machinery New York NY USA (2022). https://doi.org/10.1145/3485447.3511993
- Coustié O, Mothe J, Teste O, Baril X (2020) Meting: A robust log parser based on frequent n-gram mining, pp 84–88. https://doi.org/10.1109/ICWS49710.2020.00018
- Le VH, Zhang H (2023) Log parsing with prompt-based few-shot learning. In: 2023 IEEE/ACM 45th international conference on software engineering (ICSE) pp. 2438–2449. https://doi.org/10.1109/ICSE486 19.2023.00204
- Xiao T, Quan Z, Wang ZJ, Zhao K, Liao X (2020) Lpv: A log parser based on vectorization for offline and online log parsing. In: 2020 IEEE international conference on data mining (ICDM), pp 1346–1351. https://doi.org/10.1109/ICDM50108.2020.00175
- Zhu J, He S, Liu J, He P, Xie Q, Zheng Z, Lyu MR (2019) Tools and benchmarks for automated log parsing. In: 2019 IEEE/ACM 41st international conference on software engineering: Software engineering in practice (ICSE-SEIP), pp 121–130. https://doi.org/10.1109/ICSE-SEIP.2019.00021
- Wang X, Zhang X, Li L, He S, Zhang H, Liu Y, Zheng L, Kang Y, Lin Q, Dang Y, Rajmohan S, Zhang D (2022) Spine: a scalable log parser with feedback guidance. In: Proceedings of the 30th ACM joint european software engineering conference and symposium on the foundations of software engineering. ESEC/FSE 2022, Association for Computing Machinery New York NY USA, pp 1198–1208. https://doi.org/10.1145/3540250.3549176



173

- Du M, Li F, Zheng G, Srikumar V (2017) Deeplog: Anomaly detection and diagnosis from system logs through deep learning. In: Proceedings of the 2017 ACM SIGSAC conference on computer and communications security. CCS '17, Association for Computing Machinery New York NY USA, pp. 1285-1298. https://doi.org/10.1145/3133956.3134015
- Karlsen E, Luo X, Zincir-Heywood N, Heywood M (2024) Benchmarking large language models for log analysis security and interpretation. J Netw Syst Manag 32(3):59
- Guo H, Yuan S, Wu X (2021) Logbert: Log anomaly detection via bert. In: 2021 international joint conference on neural networks (IJCNN), pp 1–8. https://doi.org/10.1109/IJCNN52387.2021.9534113
- Le VH, Zhang H (2022) Log-based anomaly detection with deep learning: how far are we? In: Proceedings of the 44th international conference on software engineering. ICSE '22, Association for Computing Machinery, New York NY USA, pp 1356-1367. https://doi.org/10.1145/3510003.3510155
- Zhao N, Wang H, Li Z, Peng X, Wang G, Pan Z, Wu Y, Feng Z, Wen X, Zhang W, Sui K, Pei D (2021) An empirical investigation of practical log anomaly detection for online service systems. ESEC/FSE 2021. Association for Computing Machinery New York NY USA. https://doi.org/10.1145/3468264.3473933
- Zhang X, Xu Y, Lin Q, Qiao B, Zhang H, Dang Y, Xie C, Yang X, Cheng Q, Li Z, Chen J, He X, Yao R, Lou JG, Chintalapati M, Shen F, Zhang D (2019) Robust log-based anomaly detection on unstable log data. In: Proceedings of the 2019 27th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering. ESEC/FSE 2019, Association for Computing Machinery New York NY USA, pp 807–817. https://doi.org/10.1145/3338906.3338931
- Du Q, Zhao L, Xu J, Han Y, Zhang S (2021) Log-based anomaly detection with multi-head scaled dotproduct attention mechanism, pp 335-347. https://doi.org/10.1007/978-3-030-86472-9 31
- Lin CY (2004) Rouge: A package for automatic evaluation of summaries. In text summarization branches out. association for computational linguistics, Association for Computational Linguistics Barcelona Spain
- Zhou X, Peng X, Xie T, Sun J, Ji C, Liu D, Xiang Q, He C (2019) Latent error prediction and fault localization for microservice applications by learning from system trace logs. In: Proceedings of the 2019 27th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering. ESEC/FSE 2019, Association for Computing Machinery New York NY USA, pp 683-694. https://doi.org/10.1145/3338906.3338961
- He S, Lin Q, Lou JG, Zhang H, Lyu MR, Zhang D (2018) Identifying impactful service system problems via log analysis. In: Proceedings of the 2018 26th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering. ESEC/FSE 2018, Association for Computing Machinery New York NY USA, pp 60-70. https://doi.org/10.1145/3236024.3236083
- Liu Y, Yang H, Zhao P, Ma M, Wen C, Zhang H, Luo C, Lin Q, Yi C, Wang J, Zhang C, Wang P, Dang Y, Rajmohan S, Zhang D (2022) Multi-task hierarchical classification for disk failure prediction in online service systems. In: Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining. KDD '22, Association for Computing Machinery New York NY USA, pp 3438-3446. https://doi.org/10.1145/3534678.3539176
- Ma M, Liu Y, Tong Y, Li H, Zhao P, Xu Y, Zhang H, He S, Wang L, Dang Y, Rajmohan S, Lin Q (2022) An empirical investigation of missing data handling in cloud node failure prediction. In: Proceedings of the 30th ACM Joint european software engineering conference and symposium on the foundations of software engineering. ESEC/FSE 2022, Association for Computing Machinery New York NY USA, pp 1453-1464. https://doi.org/10.1145/3540250.3558946
- Luo C, Zhao P, Qiao B, Wu Y, Zhang H, Wu W, Lu W, Dang Y, Rajmohan S, Lin Q, Zhang D (2021) Ntam: Neighborhood-temporal attention model for disk failure prediction in cloud platforms. In: Proceedings of the web conference 2021. WWW '21, Association for Computing Machinery New York NY USA, pp 1181–1191. https://doi.org/10.1145/3442381.3449867
- Zhou P, Wang Y, Li Z, Wang X, Tyson G, Xie G (2020) Logsayer: Log pattern-driven cloud component anomaly diagnosis with machine learning, pp 1-10 . https://doi.org/10.1109/IWQoS49365.2020.921 2954
- Meng W, Zaiter F, Zhang Y, Liu Y, Zhang S, Tao S, Zhu Y, Han T, Zhao Y, Wang E, Zhang Y, Pei D (2023) Logsummary: Unstructured log summarization for software systems. IEEE Trans Netw Serv Manag 20(3):3803-3815. https://doi.org/10.1109/TNSM.2023.3236994
- Sui Y, Zhang Y, Sun J, Xu T, Zhang S, Li Z, Sun Y, Guo F, Shen J, Zhang Y, Pei D, Yang X, Yu L (2023) Logkg: Log failure diagnosis through knowledge graph. IEEE Trans Serv Comput16(5):3493-3507. https://doi.org/10.1109/TSC.2023.3293890
- Liu F, Wen Y, Zhang D, Jiang X, Xing X, Meng D (2019) Log2vec: A heterogeneous graph embedding based approach for detecting cyber threats within enterprise. In: Proceedings of the 2019 ACM SIGSAC conference on computer and communications security. CCS '19, Association for Computing Machinery New York NY USA, pp 1777-1794. https://doi.org/10.1145/3319535.3363224
- Locke S, Li H, Chen THP, Shang W, Liu W (2022) Logassist: Assisting log analysis through log summarization. IEEE Trans Softw Eng 48(9):3227–3241. https://doi.org/10.1109/TSE.2021.3083715



- He M, Jia T, Duan C, Cai H, Li Y, Huang G (2024) Llmelog: An approach for anomaly detection based on llm-enriched log events. In: 2024 IEEE 35th international symposium on software reliability engineering (ISSRE), IEEE, pp 132–143
- Liu Y, Tao S, Meng W, Wang J, Ma W, Chen Y, Zhao Y, Yang H, Jiang Y (2024) Interpretable online log analysis using large language models with prompt strategies. In: Proceedings of the 32nd IEEE/ACM international conference on program comprehension, pp 35–46
- OpenAI Achiam J, Adler S, Agarwal S, Ahmad L, Akkaya I, Aleman, FL, Almeida D, Altenschmidt J, Altman S, Anadkat S, Avila R, Babuschkin I, Balaji S et al (2024) GPT-4 technical report. arXiv:2303.08774
- Meng W, Zaiter F, Zhang Y, Liu Y, Zhang S, Tao S, Zhu Y, Han T, Zhao Y, Wang E, Zhang Y, Pei D (2023) Logsummary: Unstructured log summarization for software systems. IEEE Trans Netw Serv Manag 20(3):3803–3815. https://doi.org/10.1109/TNSM.2023.3236994
- THUDM (2024) Thudm/chatglm4. https://github.com/THUDM/ChatGLM4
- Bai J, Bai S, Chu Y, Cui Z, Dang K, Deng X, Fan Y, Ge W, Han Y, Huang F, Hui B, Ji L, Li M et al (2023) Qwen Technical Report . arXiv:2309.16609
- Fawcett T (2006) An introduction to roc analysis. Patt Recognit Lett 27(8):861-874
- Lin CY (2004) Rouge: A package for automatic evaluation of summaries. In text summarization branches out. association for computational linguistics. Association for Computational Linguistics Barcelona Spain
- Papineni K, Roukos S, Ward T, Zhu WJ (2002) Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the association for computational linguistics, Association for Computational Linguistics Philadelphia Pennsylvania USA
- Liang P, Bommasani R, Lee T, Tsipras D, Soylu D, Yasunaga M, Zhang Y, Narayanan D, Wu Y, Kumar A et al (2022) Holistic evaluation of language models. arXiv e-prints
- Srivastava A, Rastogi A, Rao A, Shoeb AAM, Abid A, Fisch A, Brown AR, Santoro A, Gupta A, Garriga-Alonso A, et al (2022) Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. arXiv e-prints
- Zhang L, Cai W, Liu Z, Yang Z, Dai W, Liao Y, Qin Q, Li Y, Liu X, Liu Z et al (2023) Fineval: A chinese financial domain knowledge evaluation benchmark for large language models. arXiv e-prints
- Singhal K, Azizi S, Tu T, Mahdavi SS, Wei J, Chung HW, Scales N, Tanwani A, Cole-Lewis H, Pfohl S, *et al.*: Large language models encode clinical knowledge. Nat 620(7972):172–180
- Li J, Wang X, Wu X, Zhang Z, Xu X, Fu J, Tiwari P, Wan X, Wang B (2023) Huatuo-26m,a large-scale chinese medical qa dataset. arXiv e-prints
- Miao Y, Bai Y, Li Chen HS, Dan Li Wang X, Luo Z, Sun D, Xu X, Zhang Q, Xiang C, Li, X (2023) An empirical study of netops capability of pre-trained large language models. arXiv e-prints
- Liu Y, Pei C, Xu L, Chen B, Sun M, Zhang Z, Sun Y, Zhang S, Wang K, Zhang H, Li J, Xie G, Wen X, Nie X, Ma M, Pei D (2023) Opseval: A comprehensive it operations benchmark suite for large language models. arXiv e-prints
- Silva A, Monperrus M (2024) Repairbench: Leaderboard of frontier models for program repair. arXiv preprint arXiv:2409.18952
- Jiang Z, Liu J, Chen Z, Li Y, Huang J, Huo Y, He P, Gu J, Lyu MR (2024) Lilac: Log parsing using llms with adaptive parsing cache. Proc ACM Softw Eng 1(FSE):137–160
- Sui Y, Zhang Y, Sun J, Xu T, Zhang S, Li Z, Sun Y, Guo F, Shen J, Zhang Y, Pei D, Yang X, Yu L (2023) Logkg: Log failure diagnosis through knowledge graph. IEEE Trans Serv Comput 16(5):3493–3507. https://doi.org/10.1109/TSC.2023.3293890
- Zhang W, Cheng X, Zhang Y, Yang J, Guo H, Li Z, Yin X, Guan X, Shi X, Zheng L et al (2024) Eclipse: Semantic entropy-lcs for cross-lingual industrial log parsing. arXiv preprint arXiv:2405.13548
- Liu J, Huang J, Huo Y, Jiang Z, Gu J, Chen Z, Feng C, Yan M, Lyu MR (2023) Scalable and adaptive log-based anomaly detection with expert in the loop. arXiv preprint arXiv:2306.05032
- Qi J, Huang S, Luan Z, Yang S, Fung C, Yang H, Qian D, Shang J, Xiao Z, Wu Z (2023) Loggpt: Exploring chatgpt for log-based anomaly detection. In: 2023 IEEE international conference on high performance computing & communications data science & systems smart city & dependability in sensor cloud & big data systems & application (HPCC/DSS/SmartCity/DependSys), IEEE, pp 273–280
- Shan S, Huo Y, Su Y, Li Y, Li D, Zheng Z (2024) Face it yourselves: An Ilm-based two-stage strategy to localize configuration errors via logs. In: Proceedings of the 33rd ACM SIGSOFT international symposium on software testing and analysis, pp 13–25
- Xu J, Cui Z, Zhao Y, Zhang X, He S, He P, Li L, Kang Y, Lin Q, Dang Y et al (2024) Unilog: Automatic logging via llm and in-context learning. In: Proceedings of the 46th IEEE/ACM international conference on software engineering, pp 1–12
- Karlsen E, Luo X, Zincir-Heywood N, Heywood M (2024) Benchmarking large language models for log analysis security and interpretation. J Netw Syst Manag 32(3):59



Radford A, Wu J, Child R, Luan D, Amodei D, Sutskever I (2019) Language models are unsupervised multitask learners. https://api.semanticscholar.org/CorpusID:160025533

Block J, Chen YP, Budharapu A, Anthony L, Dorr B (2023) Summary cycles: Exploring the impact of prompt engineering on large language models' interaction with interaction log information. In: Proceedings of the 4th workshop on evaluation and comparison of NLP systems, pp 85–99

Rosado T, Bernardino J (2014) An overview of openstack architecture. In: Proceedings of the 18th international database engineering & applications symposium. IDEAS '14, Association for Computing Machinery New York NY USA, pp 366–367. https://doi.org/10.1145/2628194.2628195

Jiang Z, Liu J, Huang J, Li Y, Huo Y, Gu J, Chen Z, Zhu J, Lyu MR (2024) A large-scale evaluation for log parsing techniques: How far are we? In: Proceedings of the 33rd ACM SIGSOFT international symposium on software testing and analysis, pp 223–234

OpenAI (2022) Introducing ChatGPT. https://openai.com/blog/chatgpt

Anthropic (2023). https://claude.ai/

Team G, Anil R, Borgeaud S, Alayrac JB, Yu J, Soricut R, Schalkwyk J et al (2023) Gemini: A family of highly capable multimodal models. arXiv preprint arXiv:2312.11805

Jiang AQ, Sablayrolles A, Mensch A, Bamford C, Chaplot DS, Casas D, Bressand F, Lengyel G, Lample G, Saulnier L, Lavaud LR, Lachaux MA, Stock P, Scao TL, Lavril T, Wang T, Lacroix T, Sayed WE (2023) Mistral 7B

InternLM (2023) InternLM: A multilingual language model with progressively enhanced capabilities. https://github.com/InternLM/InternLM

CodeFuse (2023). https://github.com/codefuse-ai/CodeFuse-DevOps-Model/

BAAI (2023). https://github.com/FlagAI-Open/Aquila2

Yang A, Xiao B, Wang B, Zhang B, Bian C, Yin C, Lv C, Pan D, Wang D, Yan D, Yang F, Deng F, Wang F, Liu F et al (2023) Baichuan 2: Open large-scale language models. arXiv:2309.10305

Liu Y, Tao S, Meng W, Yao F, Zhao X, Yang H (2024) Logprompt: Prompt engineering towards zero-shot and interpretable log analysis. In: Proceedings of the 2024 IEEE/ACM 46th international conference on software engineering: Companion proceedings, pp 364–365

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.



#### **Authors and Affiliations**

Tianyu Cui<sup>1</sup> · Shiyu Ma<sup>1</sup> · Ziang Chen<sup>1</sup> · Tong Xiao<sup>2</sup> · Chenyu Zhao<sup>1</sup> · Shimin Tao<sup>3</sup> · Yilun Liu<sup>3</sup> · Shenglin Zhang<sup>1,4</sup> · Duoming Lin<sup>1</sup> · Changchang Liu<sup>1</sup> · Yuzhe Cai<sup>1</sup> · Weibin Meng<sup>3</sup> · Yongqian Sun<sup>1,5</sup> · Dan Pei<sup>2</sup>

Shenglin Zhang zhangsl@nankai.edu.cn

Tianyu Cui cuitianyu@mail.nankai.edu.cn

Shiyu Ma mashiyu@mail.nankai.edu.cn

Ziang Chen 2120240792@mail.nankai.edu.cn

Tong Xiao xiaotong@tsinghua.edu.cn

Chenyu Zhao zhaochenyu@mail.nankai.edu.cn

Shimin Tao taoshimin@huawei.com

Yilun Liu liuyilun3@huawei.com

Duoming Lin 1052148783@qq.com

Changchang Liu 2113411@mail.nankai.edu.cn

Yuzhe Cai 2212113@mail.nankai.edu.cn

Weibin Meng m\_weibin@163.com

Yongqian Sun sunyongqian@nankai.edu.cn

Dan Pei peidan@tsinghua.edu.cn

- Nankai University, Tianjin, China
- <sup>2</sup> Tsinghua University, Beijing, China
- <sup>3</sup> Huawei, Beijing, China
- <sup>4</sup> Haihe Laboratory of Information Technology Application Innovation (HL-IT), Tianjin, China
- <sup>5</sup> Tianjin Key Laboratory of Software Experience and Human Computer Interaction, Tianjin, China

