

# Hue: A User-Adaptive Parser for Hybrid Logs

Junjielong Xu

siyuexi@foxmail.com

The Chinese University of Hong  
Kong, Shenzhen  
China

Yutong Cheng

222010519@link.cuhk.edu.cn

The Chinese University of Hong  
Kong, Shenzhen  
China

Qiuai Fu

fuqiuai@huawei.com

Huawei Cloud Computing  
Technologies CO., LTD.  
China

Zhijing Li

zhijingbaby@gmail.com

The Chinese University of Hong  
Kong, Shenzhen  
China

Zhouruixing Zhu

zhuzhouruixing@icloud.com

The Chinese University of Hong  
Kong, Shenzhen  
China

Yuchi Ma

mayuchi1@huawei.com

Huawei Cloud Computing  
Technologies CO., LTD.  
China

Pinjia He\*

hepinjia@cuhk.edu.cn

The Chinese University of Hong  
Kong, Shenzhen  
China

## ABSTRACT

Log parsing, which extracts log templates from semi-structured logs and produces structured logs, is the first and the most critical step in automated log analysis. While existing log parsers have achieved decent results, they suffer from two major limitations by design. First, they do not natively support hybrid logs that consist of both single-line logs and multi-line logs (e.g., Java Exception and Hadoop Counters). Second, they fall short in integrating domain knowledge in parsing, making it hard to identify ambiguous tokens in logs. This paper defines a new research problem, *hybrid log parsing*, as a superset of traditional log parsing tasks, and proposes *Hue*, the first attempt for hybrid log parsing via a user-adaptive manner. Specifically, Hue converts each log message to a sequence of special wildcards using a key casting table and determines the log types via line aggregating and pattern extracting. In addition, Hue can effectively utilize user feedback via a novel merge-reject strategy, making it possible to quickly adapt to complex and changing log templates. We evaluated Hue on three hybrid log datasets and sixteen widely-used single-line log datasets (i.e., Loghub). The results show that Hue achieves an average grouping accuracy of 0.845 on hybrid logs, which largely outperforms the best results (0.563 on average) obtained by existing parsers. Hue also exhibits SOTA performance on single-line log datasets. Furthermore, Hue has been successfully deployed in a real production environment for daily hybrid log parsing.

## CCS CONCEPTS

• **Software and its engineering** → **Software creation and management.**

## KEYWORDS

Log analysis, Software reliability

## ACM Reference Format:

Junjielong Xu, Qiuai Fu, Zhouruixing Zhu, Yutong Cheng, Zhijing Li, Yuchi Ma, and Pinjia He\*. 2023. Hue: A User-Adaptive Parser for Hybrid Logs. In *Proceedings of The 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2023)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

## 1 INTRODUCTION

In recent years, software systems, such as online services (e.g., Google Search and Bing Search) and system software (e.g., Android and Windows), have become an integral part of our daily lives, which generate extremely large amounts of software logs every day. These logs can be used in various tasks, e.g., anomaly detection [4, 15, 22, 35], root cause analysis [1, 14, 19], failure prediction [34], log compression [20, 33], and user profile construction [18]. Because of the rapid growth of the log volume, it is difficult to identify valuable information from the massive log data manually. To this end, automatic log analysis has been widely studied in recent years. The first step of automatic log analysis is log parsing, which aims at extracting log templates and converting semi-structured log messages into structured log messages for downstream tasks. Specifically, the core task of log parsing is to distinguish between *constants* and *variables* in log messages, where constants are the tokens written by developers in the logging statements (e.g., a description of a software operation) and variables are tokens that can change according to runtime environments (e.g., an IP address). Recent log parsers [3, 7, 12, 23, 28, 32] have achieved decent results on open-source log datasets (i.e., Loghub [36]). However, they still suffer from two major limitations:

*First*, existing parsers assume the incoming logs are single-line, and they try to extract the common pattern as templates line by line. However, in practice, log messages could be multi-line, such as KPI tabular echos (e.g., CPU usage), tracebacks (e.g., Java Exception), and key-value pairs (e.g., Hadoop Counters). In addition, due to the complexity of modern software and the centralized log collection

\*Corresponding author

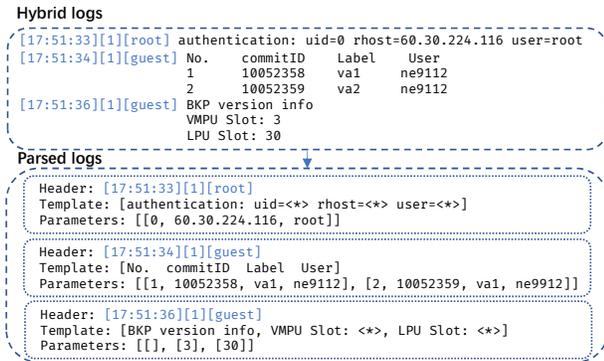


Figure 1: A simplified example of hybrid log parsing.

practice, software logs could be *hybrid*, where single-line log messages mix with multi-line log messages. For example, Fig. 1 presents an example of hybrid log parsing that contains a single-line log message and two multi-line log messages. In our opinion, hybrid log parsing is a more general research problem that aims to parse either single-line logs, multi-line logs, or their mix, while recent research mainly focuses on parsing single-line logs.

*Second*, since the logging statements of these log messages are often inaccessible (e.g., written in third-party libraries), it is difficult to tell whether a token is a constant or a variable. People could define different templates for the same logs due to the disagreement on some ambiguous tokens. For example, recent research [17, 21] reports that the widely-used parsing datasets Loghub contains labeling errors, which are typically caused by the difference in opinions. In addition, whether a log template is "correct" or not depends on the requirements of the downstream tasks. For example, in a real business scenario in Huawei Cloud, the developer-labeled log template for the console log "display ipv6" is "<\*> ipv6" in the root cause analysis task and "display <\*>" in the user profile construction task. In practice, many "correct" log templates can be summarized only with expert domain knowledge. Thus, we argue that human feedback is "the last mile" in log parsing and a log parser that can effectively and efficiently integrate human feedback is highly in demand.

In this paper, we introduce a new, general, and practical research problem, *hybrid log parsing*, which aims to parse single-line logs, multi-line logs, and their mix. To this end, we propose *Hue*, the first hybrid log parser that works in an online manner and efficiently adopts human feedback. Hue is simple yet effective. Hue adopts key casting, line aggregating, and pattern extracting to precisely parse complex hybrid logs. We further design a novel mechanism that allows users to reject a template update on potentially ambiguous tokens. We evaluate our approach on (1) three hybrid log datasets collected from both open-source software and Huawei Cloud's cloud services, and (2) sixteen widely-used single-line log datasets (i.e., Loghub). The results show that Hue achieves an average grouping accuracy (GA) of 0.845, which largely outperforms the best result obtained by existing parsers (0.563 on average). Hue also exhibits SOTA accuracy on single-line log parsing without leveraging the human feedback component, achieving the highest GA on 8 out of 16 datasets, which is the most among all compared

```

1 [17:50:30] mkdir: cannot create directory '/home/team7/results': File exists
2 [17:50:30] mkdir: cannot create directory '/home/team7/results/terasort': File exists
3 [17:50:30] mkdir: cannot create directory '/home/team7/results/bayes': File exists
4 [17:50:30] mkdir: cannot create directory '/home/team7/results/pagerank': File exists
5 [17:50:31] Traceback (most recent call last):
6 File "/HiBench3/bin/functions/load_config.py", line 713, in <module>
7   load_config(conf_root, workload_configFile, workload_folder, patching_config)
8 File "/HiBench3/bin/functions/load_config.py", line 217, in load_config
9   generate_optional_value()
10 File "/HiBench3/bin/functions/load_config.py", line 641, in generate_optional_value
11   probe_masters_slaves_hostnames()
12 File "/HiBench3/bin/functions/load_config.py", line 577, in probe_masters_slaves_hostnames
13   probe_masters_slaves_by_Yarn()
14 File "/HiBench3/bin/functions/load_config.py", line 528, in probe_masters_slaves_by_Yarn
15   assert 0, "Get workers from yarn-site.xml page failed."
16 AssertionError: Get workers from yarn-site.xml page failed.
17 [17:50:32] Parsing conf: /home/team7/HiBench3/conf/hadoop.conf
18 [17:50:32] Parsing conf: /home/team7/HiBench3/conf/hibench.conf
19 [17:50:32] Parsing conf: /home/team7/HiBench3/conf/spark.conf
20 [17:50:32] Parsing conf: /home/team7/HiBench3/conf/workloads/micro/terasort.conf
21 [17:50:34] display configuration commit list
22 -----
23 No.      CommitID      Label      User      TimeStamp
24 -----
25 1        1001225122    td-begin_backroll    root      2022-11-20 17:50:32
26 2        1001225121    -          root      2022-11-20 17:50:31
27 -----

```

Figure 2: An example of hybrid logs collected from an industrial log management platform.

parsers. Hue has been successfully deployed in the cloud services in Huawei Cloud to parse daily hybrid logs.

This paper makes the following main contributions:

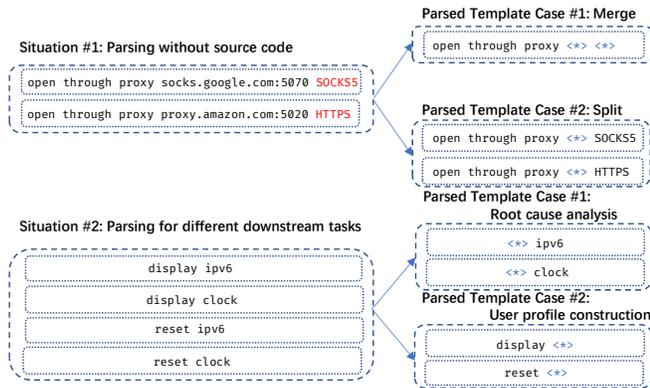
- It introduces a new, general, and practical research problem, hybrid log parsing, which is a superset of the existing single-line log parsing problem.
- It proposes Hue, the first log parser that natively supports hybrid log parsing and has a human feedback integration mechanism that largely reduces unnecessary queries.
- It presents the evaluation of Hue using three hybrid log datasets and sixteen single-line log datasets, demonstrating that Hue achieves SOTA accuracy and efficiency on both hybrid and single-line logs.
- It releases two hybrid log datasets collected from open-source software and cloud system for this research direction.

## 2 MOTIVATION

This section intends to further explain the two major limitations of existing parsers that motivates our work. The following is mainly inspired by and summarized from our collaboration experience with engineers in Huawei Cloud.

### 2.1 Neglected Hybrid Logs

Hybrid logs are a common type of data in the software industry. As shown in Fig. 2, they are usually generated from cloud platforms and can contain a mixed combination of single-line and multi-line log messages. The hybrid nature of the logs is often caused by permission restrictions between different service departments within an organization. In particular, IT operators may not be able to access logs in all components directly and must instead use a log aggregation system to gather all the service outputs for centralized management. This means that hybrid logs may contain a variety of log types, including single-line logs such as component event logs and multi-line logs such as console echoes, exception traceback, and even tabular system key performance indicators (KPIs). Hybrid logs are not only found in closed-source software, but they are also commonly encountered in open-source components such as Hadoop, Spark, and MySQL. This highlights the widespread availability of hybrid logs in both closed-source and open-source software.



**Figure 3: Two typical situations that difficult for the parser to determine the template automatically.**

Hybrid logs contain valuable information that can be utilized in many downstream tasks, including IT operators using the traceback for failure prediction and root cause analysis, QA engineers selectively testing APIs through KPI tables, and commercialization departments constructing user profiles from console echoes to target potential user needs. However, the structural variability of hybrid logs presents a challenge for current log parsers, leading to them ignoring these logs during training and deployment. If single-line parsers are used in pipelines, they tend to split a multi-line log message into multiple lines of text and regard each line as a separate single-line log message, leading to incorrect parsing results. Even if IT operators try to use special regex to flatten hybrid logs, tabular messages may still not be parsed correctly because of the uncertainty in the number of table rows. To address these issues, this paper proposes key casting (Sec. 4.2) and line aggregating techniques (Sec. 4.3). The limitations of existing work are further discussed in our experiments (Sec. 5.2).

## 2.2 Ambiguous Template Tokens

Recent research found that different people could generate different log templates because of ambiguous tokens [17, 21]. In Huawei Cloud, there are two typical situations that confuse parsers. (1) Whether to merge two candidate templates by turning a constant into a variable. For example, in Fig. 3, Situation 1, it is debatable whether the two logs belong to the same event (and should be merged). (2) A log message could have multiple "correct" templates when it was used by different downstream tasks. For example, in Fig. 3, Situation 2, for user profiling modeling, the analyst cares about the instructions and actions the user performed, whereas for root cause analysis, the analyst cares about the object being manipulated. The difference in their preferences leads to different "correct" log templates.

To improve the compatibility and generalizability of the log templates generated in log parsing, we introduce an efficient manual feedback mechanism (Sec. 4.5) that works seamlessly within a parse tree and effectively improves parsing accuracy with little manual involvement. The benefits of this mechanism are further explained in our experiments (Sec. 5.4).

## 3 PRELIMINARY

This section intends to introduce the main concepts in hybrid log parsing, especially our definition and the characteristics of *hybrid log* and the input and output of *hybrid log parsing*.

### 3.1 Hybrid Logs

Hybrid logs could be single-line logs, multi-line logs, or their mix. Based on our experience with hybrid log analysis in Huawei Cloud, in this paper, we mainly consider three kinds of logs in hybrid logs: *event logs* (single-line), *table logs* (multi-line), and *text logs* (multi-line). For example, in Fig. 2, line 1-4 and line 17-20 are event logs, line 5-16 is a text log, and line 21-27 is a table log.

**3.1.1 Event Log.** Event logs are single-line log messages. An event log records an operation or status of a service or component and contains a log header and message content. To align with the definition in existing single-line log parsing research [13, 17], an event log's template consists of constants in the log message and wildcards that indicate variables.

**3.1.2 Table Log.** Table logs are multi-line logs that contain a table header, multiple lines of parameters, and potentially table lines. Table logs with the same log template (1) have the same number of table columns and the token data type in the same column is the same (e.g., the timestamps in the last column in the table log in Fig. 2); (2) might have different numbers of rows. Typical table logs contain various system metrics and these logs are often generated by the performance testing component in the cloud and the echoes from the user shell. Since the header of the table logs consists of the key information (i.e., column number and column types) that distinguishes between different tables, the template of a table log is the table header, while the table content is regarded as parameters in the parsed log message (as illustrated in Fig. 1). Therefore the goal of log parsing on table logs is to extract its table header and transform its table content into parameters.

**3.1.3 Text Log.** Text logs are multi-line logs that do not fall into the "table log" category. Text logs record the detailed information in addition to the system event in plain text, such as traceback call stack in a program or database content in the form of key-value pairs. Text logs are widely available in many components such as Hadoop, Spark, and MySQL. However, to the best of our knowledge, only a little log analysis research utilized this kind of log and we think one of the main reasons is that existing log parsers cannot accurately parse text logs. We define the template of a text log as tokens in several lines that remain constant, such as the traceback error type or the keys in the message. Thus, the goal of log parsing on text logs is to extract these lines from the message.

### 3.2 Hybrid Log Parsing

The goal of hybrid log parsing is to extract log templates from log messages and generated structured logs accordingly. In our opinion, hybrid log parsing, which focuses on single-line logs, multi-line logs, and their mix, is a *superset* of existing log parsing research that targets single-line logs. The output of this task includes (1) a meta-file that contains the template of each log message and its

location in the original log file, and (2) files with structured variable lists (e.g., headers and parameters) in each log type.

In this paper, we view hybrid log parsing as a classification problem. Given a raw log message  $L$  consisting of  $n$  tokens across  $l$  lines, denoted as  $[t_1, t_2, \dots, t_n]$ , a hybrid log parser should predict the type  $T$  of log message  $L$ , the log template group  $g_k^T \in G^T = \{g_1^T, g_2^T, \dots, g_N^T\}$  it belongs to among the total  $N$  templates in type  $T$ , and whether a token  $t_i$  is a constant ( $y_i = 1$ ) or a variable ( $y_i = 0$ ). Finally, the log parser must output its predictions in a structured format. Hybrid log parsing is challenging because it is non-trivial to automatically distinguish between different types of logs when they are mixed and identify the constants and variables in them.

## 4 APPROACH

### 4.1 Overview

This section introduces Hue, an unsupervised online log parser that utilizes a heuristic parse tree and an efficient human feedback integration mechanism. The overview of Hue is presented in Fig. 4, which contains four main components: key casting, line aggregating, pattern extracting, and online updating.

Specifically, Hue streamly reads the log data line by line, splits each line into tokens by spaces, and sends the raw *token sequence* to main components for parsing: First, *key casting* transforms part of the tokens into *keys*, i.e., special wildcards indicating representative token data type. For instance, a raw token sequence with  $n$  tokens  $S_{raw} = [t_1, t_2, t_3, \dots, t_n]$  can be casted to a new token sequence  $S = [t_1, \mathcal{K}_2, \mathcal{K}_3, \dots, t_n]$ , where each  $\mathcal{K}$  represents a specific key. When all lines in a log are fully casted, *line aggregating* combines adjacent token sequences into *blocks* (i.e., lists of adjacent token sequence) via heuristic rules. For example, a log with  $l$  lines  $L = [S_1, S_2, \dots, S_l]$  can be aggregated to  $L' = [B_1, B_2, \dots, B_k] = [[S_1], [S_2, S_3], \dots, [S_l]]$ , which contains  $k$  blocks. Then, *pattern extracting* utilizes aggregated blocks to obtain log templates. Finally, *online parsing* revises the current log template in an online manner with the guidance of an optional human feedback mechanism.

Note that hybrid logs could be single-line logs, multi-line logs, or their mix, thus Hue, which is the first hybrid log parser, can be utilized to parse single-line logs as existing parsers.

### 4.2 Key Casting

Preprocessing has been widely adopted by existing parsers. A typical preprocessing process identifies common variables (e.g., IP address) by regular expressions, and then either remove the identified tokens or replaces them with the wildcard "<\*>" [13]. However, these preprocessing strategies are not effective when dealing with hybrid logs, which may contain table logs and text logs. Table logs have lines with only variables, which would be all removed or transformed into the same wildcard "<\*>" by existing preprocessing strategies, leading to significant inaccuracy. In addition, these strategies might exacerbate the difficulty of distinguishing text logs from table logs. To address this problem, it is possible to manually configure a list of delimiters for each system separately. However, this is impractical for hybrid logs because they are often collected from different system components or even different systems.

Hue adopts a novel preprocessing strategy called *key casting*. The main idea is to use different wildcards (e.g., "<\*>" and "<\*>")

to better encodes prior knowledge in preprocessing. Hue split a line by spaces and transforms tokens that have been commonly used as variables (e.g., IP addresses, file paths, boolean values) into the corresponding wildcards called "keys" according to a general casting table. For example, in Fig. 4, "authentication: uid=0 rhost=60.30.224.116 user=root" would be transformed into "authentication: uid=<\*> rhost=<\*> user=root". Key casting is simple yet effective for hybrid log parsing, and we believe it could also benefit existing log parsing approaches as a general preprocessing strategy. In our implementation, we construct a general key casting table for all datasets, which contains seven kinds of keys, as illustrated in Fig. 4. Users can also easily customize their own key casting table by changing one line of code.

### 4.3 Line Aggregating

To correctly parse hybrid logs, identifying log type is critical. Intuitively, using line numbers to identify event logs is relatively easy; but distinguishing between table logs and text logs is challenging because both of them have a multi-line structure. To overcome this issue, we propose *line aggregating* to differentiate between table logs and text logs. Specifically, we observe that adjacent lines in table logs tend to exhibit higher similarity, whereas in text logs, although adjacent lines are less similar, they often share the same indentation. Thus, Hue aggregates adjacent lines in multi-line logs and identifies log types based on adjacent lines' (1) sequence similarity and (2) indent number.

The workflow of line aggregating is shown in Algo. 1. Hue first determines whether the current log is an event log by checking if it only consists of single line (line 2). If the log has multiple lines, Hue sequentially aggregates adjacent lines into several *blocks* (Sec. 4.1) if sequence similarity exceeds the aggregating similarity threshold  $\epsilon_a$  (i.e., a hyperparameter determined before parsing) or they have the same indentation (line 4-13). Then, Hue utilizes a *type counter*  $C$  to keep track of the predominant reason in the aggregation process (line 14): if  $C < 0$ , sequence similarity is the primary reason for aggregating, suggesting that the log should be a table log; if  $C \geq 0$ , indentation number is the primary reason for aggregating, suggesting that the log should be a text log. Finally, Hue utilizes the aggregated logs for subsequent pattern extraction. To enhance understandability, we also provide an example of line aggregating for text logs in Fig. 5. Specifically, similarity  $Sim(S_1, S_2)$  represents the proportion of identical tokens at the same positions in two sequences (Eq. 1), where  $N$  is their average token length, and  $Equ(t_1, t_2)$  measures the equality of two tokens (Eq. 2).

$$Sim(S_1, S_2) = \frac{\sum_{i=1}^N Equ(S_1[i], S_2[i])}{N} \quad (1)$$

$$Equ(t_1, t_2) = \begin{cases} 1 & \text{if } t_1 = t_2 \\ 0 & \text{if } t_1 \neq t_2 \end{cases} \quad (2)$$

Notably, there are a few details in implementing line aggregating: (1) Lines consisting of only delimiters or null characters should be skipped. (2) Based on our observation in both open-source and industrial software, hybrid logs are often preceded by log headers, which can serve as delimiters between logs.

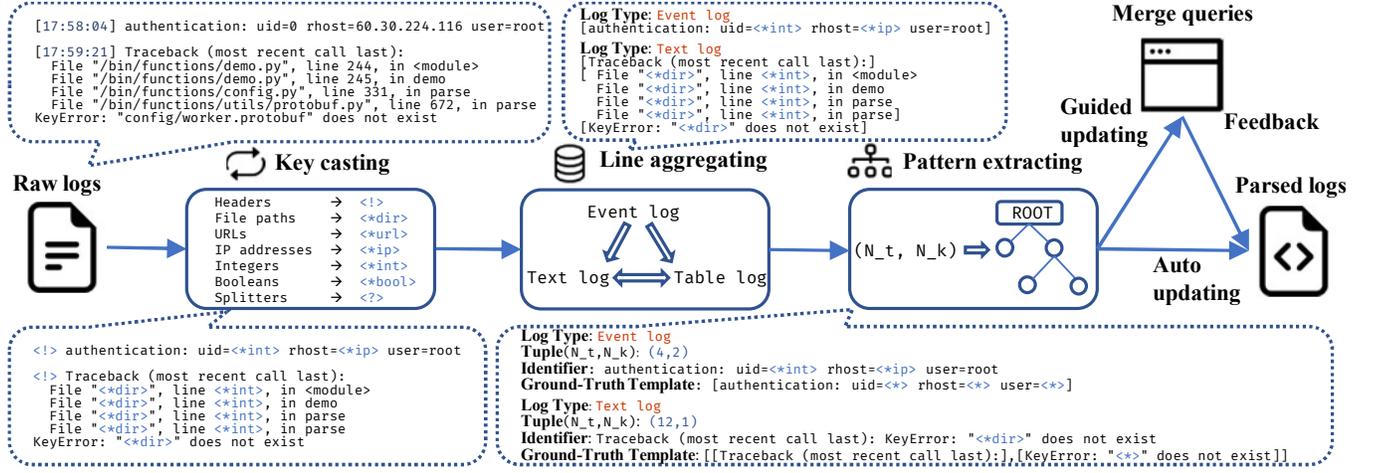


Figure 4: The overview of Hue with two examples. Users can customize the key casting table and provide feedback to optimize Hue for more complex data types and usage scenarios based on the target data characteristics and their own domain expertise.

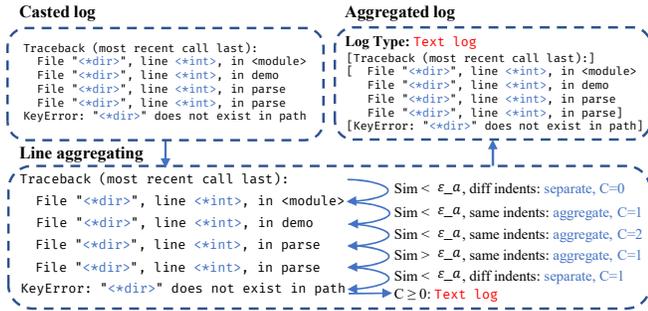


Figure 5: An example of line aggregating. (assume  $\epsilon_a = 0.9$ )

#### Algorithm 1 Line Aggregation & Type Determination

**Input:** Sequence queue of current message,  $Q = \{S_1, S_2, \dots, S_n\}$   
 Aggregating similarity threshold  $\epsilon_a$

- 1: Log type  $\mathcal{T} = \text{"EVENT"}$ , Block queue  $Q_B = \{\}$ , Last block  $B_{last} = S_1$
- 2: **if**  $|Q| > 1$  **then**
- 3:   Type counter  $C = 0$
- 4:   **for**  $S_i \in Q$  **do**
- 5:      $S_{B_{last}}$  = the common sequence of all  $S \in B_{last}$
- 6:     **if**  $\text{Sim}(S_i, S_{B_{last}}) \geq \epsilon_a$  **then**
- 7:        $C \leftarrow C - 1$ ,  $B_{last} \leftarrow B_{last} \cup S_i$
- 8:     **else if**  $S_i$  and  $S_{B_{last}}$  have the same indents **then**
- 9:        $C \leftarrow C + 1$ ,  $B_{last} \leftarrow B_{last} \cup S_i$
- 10:    **else**
- 11:      $Q_B \leftarrow Q_B \cup B_{last}$ ,  $B_{last} = \{S_i\}$
- 12:    **end if**
- 13:   **end for**
- 14:    $\mathcal{T} = \text{"TABLE"}$  if  $C < 0$ , else:  $\mathcal{T} = \text{"TEXT"}$
- 15: **end if**
- 16:  $Q_B \leftarrow Q_B \cup B_{last}$

**Output:**  $Q_B, \mathcal{T}$

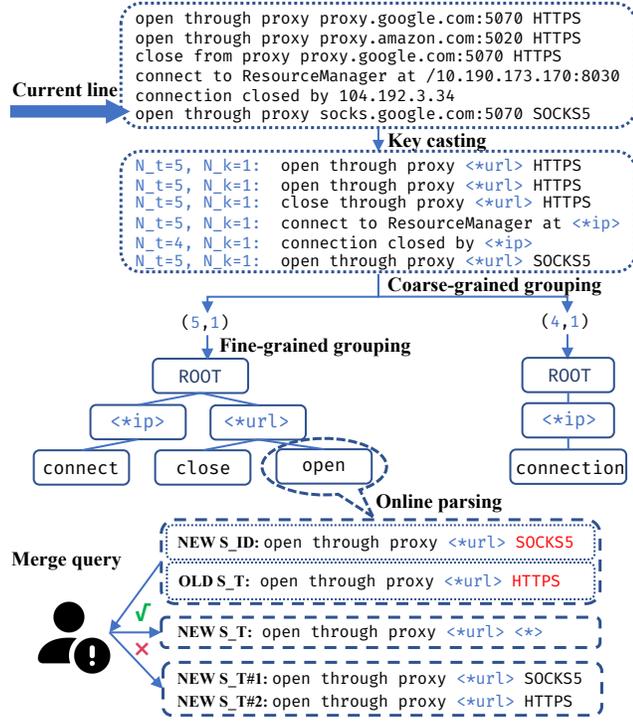
## 4.4 Pattern Extracting

To extract log templates, log parsers often first group logs with the same template into the same group [12, 28, 32]. Traditional log parsers often use complete log messages (all tokens) for clustering, which is inefficient in hybrid log parsing (shown in Fig. 9) because a hybrid log message can be very long after flattened. In addition, the line number of multi-line logs from the same template group can vary greatly, leading to low accuracy in message-level clustering.

To tackle these challenges, Hue proposed *pattern extracting*, which utilizes the block queue  $Q_B$  and log type  $T$  obtained from line aggregating to group logs with the same templates. As shown in Fig. 4, Hue uses a special sequence of each log, the "log identifier" ( $S_{ID}$ ), which may contain a log template as defined in Sec. 2, as a manual feature. Hue assumes that logs from the same template have similar  $S_{ID}$ , and uses this information to group logs. Specifically,  $S_{ID}$  is (1) the entire sequence for an event log, (2) the concatenated common sequences from the first and last blocks for a text log, and (3) the sequence of the first block with a capacity of one before the first block with a capacity greater than one for a table log. The choice of the identifier is based on an assumption that text logs often have template tokens at the beginning or end (e.g., Error Type), and that table headers are typically distinct and should be aggregated with the table content. 4 as shown in Fig.6.

**4.4.1 Coarse-grained Grouping.** To group logs, Hue counts the number of tokens and keys in each log's identifier  $S_{ID}$ , creating a tuple  $(N_t, N_k)$ . The logs are then assigned to the group with the same tuple. This process is based on two assumptions: (1) identifiers from the same template have the same token number; (2) identifiers from the same template have the same key number. The first assumption is commonly used and has been verified in previous research [12, 23], where the identifiers of event logs are themselves.

**4.4.2 Fine-grained Grouping.** To further group the logs, Hue uses a multinomial search tree with a unique search rule for each tuple  $(N_t, N_k)$ . First, the tree forks based on the keys in the  $N_k$  layer, and then forks based on non-key prefix tokens in the next  $N_p$



**Figure 6: An example of pattern extracting and guided online parsing for event logs. (assume maximum tree depth is 2)**

( $N_p = N_t - N_k$ ) layer. The tree is initialized with only one root node. Hue creates a new node each time when a new identifier  $S_{ID}$  comes and no corresponding node exists at the current layer. This forking strategy is based on the assumption that template tokens are often at the beginning and variable tokens are often at the end of logs. To further improve the accuracy of parsing, Hue allows setting limits on tree depth and identifier length. This helps prevent too long or too short logs from affecting the parsing accuracy by creating incorrect parsing rules. Identifiers not meeting the length criteria will be assigned to a designated group for similarity matching, rather than going through the tree.

## 4.5 Online Parsing

Upon a  $S_{ID}$  reaching a leaf node, Hue performs *online parsing*, which involves obtaining current log's template and updating the outdated template in each log group. It assigns the current log to a specific log group, obtains the group template  $S_T$  (i.e., the common sequence of logs'  $S_{ID}$  in the group) as the current log's template, and dynamically updates  $S_T$  with new incoming  $S_{ID}$ . Previous studies [12, 23, 28] commonly employed automatic template updates and maintenance, which may be ineffective for parsing logs with ambiguous tokens, as discussed in Sec. 2.2. Therefore, Hue introduces a *guided updating* mechanism to enhance parsing performance for log files with ambiguous logs, while retaining the classic *auto updating* mechanism for regular log files. The following sections provide a detailed explanation of the two updating methods.

**4.5.1 Auto Updating.** Hue calculates the similarity (Eq. 1) between  $S_{ID}$  and all groups' templates within current node. If a template  $S_T$  has the highest similarity score that exceeds the merging similarity threshold  $\epsilon_m$  (i.e., a hyperparameter determined before parsing), Hue will update the template to be the common sequence of  $S_{ID}$  and the original template  $S_T$  (i.e.,  $S_T \leftarrow S_{ID} \cap S_T$ ). If no such templates exist, Hue creates a new template and initializes it as  $S_{ID}$ . By continuously updating and outputting log templates based on the newly added  $S_{ID}$ , Hue achieves log parsing in an online manner.

**4.5.2 Guided Updating.** Hue still calculates the similarity (Eq. 1) between  $S_{ID}$  and all templates in node. However, instead of solely relying on  $\epsilon_m$  to perform updating, Hue introduces a "merge-reject" strategy that allows users to reject a template merge on a potentially ambiguous token in candidate templates. Specifically, Hue first conducts a check: *whether tokens requiring updating in  $S_T$  contain non-key tokens?* If yes, it triggers a merge query, i.e., output current  $S_T$  and  $S_{ID}$ , and prompts the user to decide whether to reject the automatic merging and updating process. If rejected, Hue creates a new group under the same leaf node, initializing it with  $S_{ID}$  as the template. Otherwise, the merging process continues, and  $S_T$  is updated with the common sequence of tokens between  $S_T$  and  $S_{ID}$ . The two cases (i.e., trigger a merge query or not) are illustrated with examples in Fig. 7.

The idea of developing a human feedback mechanism was inspired by a recent parser, SPINE [32]. However, Hue provides a different and more efficient mechanism. Specifically, Hue largely reduces the frequency of triggered queries by skipping the updates of keys. This is because the keys are casted parameter tokens in the preprocessing and the ambiguous tokens are non-key tokens. In addition, Hue does not involve numerical computation to determine whether a merging process should trigger a query in guided mode, making it efficient in development scenarios.

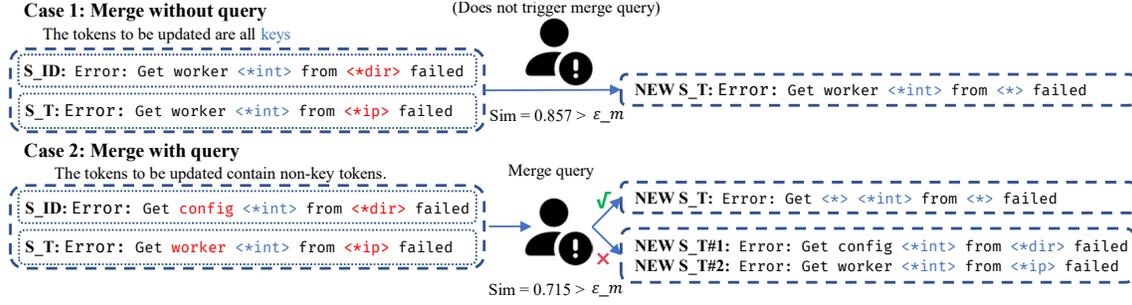
## 5 EXPERIMENTS

We conducted extensive experiments on public datasets and industrial datasets to answer the following research questions:

- **RQ1:** How accurate and efficient is Hue on hybrid logs?
- **RQ2:** How accurate and efficient is Hue on single-line logs?
- **RQ3:** How effective is Hue's human feedback integration mechanism?

### 5.1 Experiment Setup

**5.1.1 Implementation.** All of the experiments were performed on a virtual machine with 128 Intel(R) Xeon(R) Platinum 8375C CPU @ 2.90GHz processors and 94GB RAM on Ubuntu 20.04.5 LTS. We implemented Hue in Python 3.9.12 with the same key casting table shown in Fig. 4. Additionally, we implemented specific hyperparameter configuration files to manage all Hue's hyperparameters for each log source. Specifically, the hyperparameters include (1) the maximum tree depth, (2) the maximum and minimum identifiers length, (3) the aggregating similarity threshold  $\epsilon_a$ , and (4) the merging similarity threshold  $\epsilon_m$ . Similar to previous works [3, 12, 32], hyperparameters are only related to the log file source and remain unchanged during the parsing process after being determined by



**Figure 7: Two cases of guided updating mode. The updates of keys do not need human guidance since they are typically considered as non-ambiguous tokens. (assume  $\epsilon_m = 0.7$ )**

grid search, a common hyperparameter tuning method, on small-scale homogeneous log data (*i.e.*, 100 random samples for each dataset).

Notably, in experiments, we observed that setting the  $\epsilon_m$  under the guided updating mode (denoted as  $\epsilon_{mg}$ ) to be approximately 0.2 lower than  $\epsilon_m$  under the auto updating mode (denoted as  $\epsilon_{ma}$ ) will result in better parsing performance for most log sources. To avoid frequent editing of configuration files when switching updating modes, we introduced an elastic variable  $\epsilon_e = \epsilon_{ma} - \epsilon_{mg}$  to optimize the performance in both modes. Specifically,  $\epsilon_e$  is typically set to 0.2 by default. Only in a small amount of log sources,  $\epsilon_e$  need to be changed to further enhance parsing performance.

**5.1.2 Datasets.** Our experiments are conducted on 19 log datasets, including 16 single-line log datasets from LogPAI [36] and 3 hybrid log datasets collected by us. The LogPAI datasets cover a variety of log types, including OS logs, app logs, and microserver logs, with 2k messages each and 5 to 200 templates. The hybrid datasets consist of runtime logs generated by a cloud system benchmark (HiBench), multi-source logs collected by a cloud testing system (CTS), and multi-terminal logs collected in PaaS product in Huawei Cloud (Paas). We manually labeled all hybrid log datasets and released the first two datasets for replication. Table. 1 presents the statistics of these datasets. All of these datasets include ground truth labels.

**Table 1: Dataset statistics of hybrid log datasets.**

Num. (Event/Table/Text)	HiBench	CTS	PaaS
Message Num.	1879/2057/64	260/17/101	386/36/255
Template Num.	92/7/18	93/7/43	129/6/93

**5.1.3 Metrics.** Our evaluation uses two commonly used metrics: Grouping Accuracy [2, 17, 36] and Template F1-score [17]. The former measures the proportion of correctly parsed log messages, while the latter measures the F1 score of the generated templates from the parser. We use these metrics to conduct fair comparisons at both the message and template levels. Additionally, we compare the execution time of different parsers to evaluate their efficiency.

**5.1.4 Baselines.** We selected AEL [7], LenMa [28], Spell [3], IPLoM [23], Drain [12], and SPINE [32], the top-performing unsupervised

**Table 2: Grouping accuracy on hybrid log datasets.**

Metrics	AEL	LenMa	Spell	IPLoM	Drain	Hue
HiBench	0.308	0.238	0.360	0.300	0.442	<b>0.932*</b>
CTS	0.620	0.432	0.571	0.513	0.746	0.848*
PaaS	0.178	0.300	0.130	0.345	0.502	0.754*
Average	0.369	0.323	0.354	0.386	0.563	0.845*

log parsers, as our baselines for comparison. Since Hue is the first attempt for parsing hybrid logs, there is no other baseline for comparison on parsing multi-line logs. To demonstrate Hue’s effectiveness on hybrid logs, we collect and flatten multiple lines into individual log entries with line delimiters “/n”, which is a typical solution for managing multi-line logs in real development scenarios [5, 6, 26, 27] (including in Huawei Cloud) before inputting hybrid logs into existing parsers. Note that since SPINE is not open-source, we cannot reproduce their results. To improve the scientific rigor of our research, we compare Hue with SPINE in part of our experiments by using results from its original paper under the same experimental settings, *i.e.*, grouping accuracy on single-line logs and feedback results on the Linux dataset. All baselines are implementations provided by LogPAI or their original repositories with all parameters set to their optimal configuration.

## 5.2 RQ1: How Accurate and Efficient is Hue on Hybrid Logs?

For RQ1, we use grouping accuracy, template F1-score, and execution time as evaluation metrics. We conducted experiments on 3 hybrid log datasets. Note that the comparison is not *apple to apple* because as to our knowledge, Hue is the first parser attempting to parse hybrid logs (including multi-line logs). Therefore, we cannot find other more suitable baselines for comparison. Our purpose is not to claim this is a weakness of existing single-line parsers because they were not designed for hybrid log parsing, but rather to highlight the importance of designing a new parser for parsing hybrid logs. In addition, the experiments intended to show that it is non-trivial to adapt single-line parsers to hybrid logs. To be fair, we use the same set of regexes for all parsers including Hue.

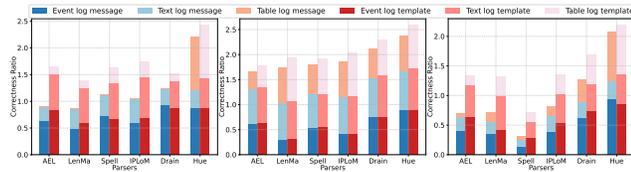
**5.2.1 Accuracy.** The results are shown in Table. 2 and Table. 3. The best results in each dataset are marked with a star symbol, and the

**Table 3: Template F1-score on hybrid log datasets.**

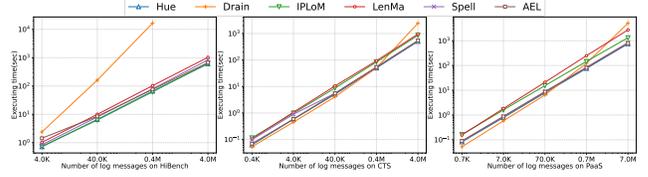
Metrics	AEL	LenMa	Spell	IPLoM	Drain	Hue
HiBench	0.083	0.590	0.641	0.687	0.753	<b>0.836*</b>
CTS	0.662	0.569	0.618	0.645	0.784	<b>0.839*</b>
PaaS	0.326	0.388	0.292	0.450	0.540	<b>0.801*</b>
Average	0.357	0.516	0.517	0.594	0.692	<b>0.825*</b>

results with grouping accuracy greater than 0.9 and F1-score greater than 0.8 are emphasized in bold font. Hue demonstrates state-of-the-art accuracy and F1-score on all three hybrid log datasets, outperforming the best traditional parser, Drain. In comparison, Drain achieves an average grouping accuracy of 0.563 and an average F1-score of 0.692, while Hue achieves the highest average grouping accuracy of 0.845 and the highest average F1-score of 0.825. This superior performance can be attributed to the design of Hue, which is specifically built to handle the unique challenges posed by hybrid log parsing. Unlike traditional parsers, Hue is able to effectively parse logs that span multiple lines and retain the structure between lines, whereas traditional parsers fall short because they were not designed to natively support multi-line log parsing. As mentioned in Sec. 2.1, the structure between lines is lost after flattening logs across lines. For example, in table logs, different instances of the same template may have different line numbers, which can result in significant differences in the message lengths after flattening. This, in turn, results in far more templates than ground truth, making it difficult or even impossible for downstream tasks to use the logs for vectorized representation.

In particular, we find that Hue has a significant advantage in parsing event logs and table logs. Fig. 8 illustrates the comparative advantages, *i.e.*, calculating the ratio of correctly parsed message number or template number dividing the total message number or template number on each log type when they are parameterized to maximize template F1-score. For example, on HiBench, Hue demonstrates promising performance in parsing table logs, while other log parsers struggle to parse table logs. In addition, on CTS, Hue exhibits significant advantages in parsing event logs. These results demonstrate the necessity of designing a parser specifically for hybrid logs, and migrating existing parsers for hybrid log parsing is not trivial.

**Figure 8: The ratio of correctly parsed message/template numbers in each type. The upper bound of each stacked bar is 3.0.**

**5.2.2 Efficiency.** We specially collected a large quantity of log data on these three hybrid logs for efficiency evaluation. Specifically, we compare average time consumption after running five times in a scaling-up scenario. The results of the experiments are displayed in

**Figure 9: Execution time on hybrid logs.****Table 4: Grouping accuracy on LogPAI’s loghub datasets.**

Dataset	AEL	LenMa	Spell	IPLoM	Drain	SPINE	Hue
HDFS	<b>0.998</b>	<b>0.998</b>	<b>1*</b>	<b>1*</b>	<b>0.998</b>	<b>0.998</b>	<b>0.998</b>
Spark	<b>0.905</b>	0.884	<b>0.905</b>	<b>0.920</b>	<b>0.920</b>	<b>0.925</b>	<b>0.942*</b>
BGL	0.758	0.690	0.787	<b>0.939</b>	<b>0.963*</b>	<b>0.948</b>	0.849
Windows	0.690	0.566	<b>0.989</b>	0.567	<b>0.997*</b>	<b>0.990</b>	<b>0.990</b>
Linux	0.690	0.701	0.605	0.672	0.690	0.676	0.749*
Android	0.682	0.880	<b>0.919</b>	0.712	<b>0.911</b>	<b>0.932*</b>	0.826
Mac	0.764	0.698	0.757	0.673	0.787	0.789	<b>0.901*</b>
Hadoop	0.538	0.885	0.778	<b>0.954</b>	<b>0.948</b>	<b>0.946</b>	<b>0.966*</b>
HealthApp	0.568	0.174	0.639	0.822	0.780	<b>0.988*</b>	<b>0.903</b>
OpenSSH	0.538	<b>0.925*</b>	0.554	0.802	0.788	0.681	0.804
Thunderb.	<b>0.941</b>	<b>0.943</b>	0.844	0.663	<b>0.955</b>	<b>0.964*</b>	<b>0.962</b>
Proxifier	0.518	0.508	0.527	0.515	0.527	<b>0.967</b>	<b>1*</b>
Apache	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>
HPC	<b>0.903</b>	0.830	0.654	0.824	0.887	0.871	<b>0.951*</b>
Zookeeper	<b>0.921</b>	0.841	<b>0.964</b>	<b>0.962</b>	<b>0.967</b>	<b>0.989*</b>	<b>0.987</b>
OpenStack	0.758	0.743	0.764	0.871	0.733	0.757	<b>0.993*</b>
Average	0.754	0.721	0.751	0.777	0.865	<b>0.901</b>	<b>0.927*</b>

Fig. 9. We observe that Hue obtains SOTA results and its execution time increases linearly with the number of logs. Some multi-line logs, after being flattened, can be very long and cause sequence clustering-based parsers to fail. Although Drain performs the best effectiveness among all baselines, it still lacks of efficiency on hybrid logs. Specifically, Drain times out when the number of messages reached 4.0 million. The reason might be that Drain relies on direct comparison to group extra-long logs, which is highly inefficient on flattened multi-line logs.

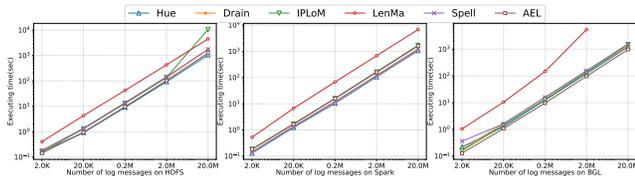
### 5.3 RQ2: How Accurate and Efficient is Hue on Single-line Logs?

For RQ2, we use grouping accuracy, template F1-score, and execution time for evaluation. We compared Hue (auto-updating) with 6 SOTA log parsers on 16 open-source single-line log datasets.

**5.3.1 Accuracy.** The results are presented in Table 4 and Table 5. Hue achieves the SOTA results with the highest average values in both grouping accuracy and template F1-score. SPINE [32] is the best among existing parsers on single-line logs. Compared with SPINE, Hue performs better than at the message level (Grouping accuracy) on 8 datasets and is comparable on 4 datasets. Hue also achieves SOTA results at the template level for 9 datasets. To sum up, Hue achieves the highest average values of 0.927 in grouping accuracy and 0.870 in template F1-score, surpassing the best traditional parser results of 0.901 and 0.749, respectively.

**Table 5: Template F1-score on LogPAI’s loghub datasets.**

Dataset	AEL	LenMa	Spell	IPLoM	Drain	Hue
HDFS	<b>0.998</b>	<b>0.867</b>	<b>1*</b>	<b>1*</b>	<b>0.867</b>	<b>0.867</b>
Spark	0.511	0.373	0.510	0.769	0.780	<b>0.870*</b>
BGL	<b>0.819</b>	0.412	0.487	<b>0.832*</b>	<b>0.830</b>	0.700
Windows	<b>0.833</b>	0.752	0.777	0.774	<b>0.932*</b>	<b>0.811</b>
Linux	<b>0.818</b>	<b>0.944*</b>	<b>0.802</b>	0.796	<b>0.930</b>	<b>0.926</b>
Android	0.726	<b>0.838</b>	<b>0.867*</b>	0.703	<b>0.825</b>	<b>0.834</b>
Mac	0.759	0.655	0.634	0.743	0.797	<b>0.885*</b>
Hadoop	0.702	0.672	0.604	<b>0.863</b>	<b>0.820</b>	<b>0.880*</b>
HealthApp	0.306	0.114	0.367	0.463	0.351	<b>0.960*</b>
OpenSSH	0.571	<b>0.824*</b>	0.630	0.667	<b>0.824*</b>	0.767
Thunderb.	0.695	0.748	0.638	0.784*	0.773	0.779
Proxifier	0.400	0.154	0.134	0.600	0.534	<b>1*</b>
Apache	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>	<b>1*</b>
HPC	0.750	0.297	0.739	0.711	0.753	<b>0.872*</b>
Zookeeper	0.718	0.755	0.729	0.744	<b>0.854*</b>	<b>0.854*</b>
OpenStack	0.657	0.244	0.089	0.693	0.117	<b>0.978*</b>
Average	0.696	0.603	0.629	0.759	0.749	<b>0.870*</b>

**Figure 10: Execution time on single-line logs.**

We also observe that Hue’s accuracy and F1-score are lower on certain datasets (e.g., Linux). We think this is caused by some ambiguous tokens in the logs. For instance, "usbcore: registered new driver *usbfs*" and "usbcore: registered new driver *hub*" result in lower accuracy on this dataset. To address this issue, Hue provides the human feedback integration mechanism (guided updating), which will be discussed in the following section.

**5.3.2 Efficiency.** We assessed the efficiency the parsers on three major datasets, i.e., using HDFS, Spark, and BGL for executing time comparisons [32]. We also compare average time consumption after running five times in a scaling-up scenario. The results are shown in Fig. 10. The results indicate that Hue is comparable with the best baseline in terms of processing speed for single-line log data. However, Hue is slightly slower than the fastest baseline in some cases, which we argue is acceptable.

#### 5.4 RQ3: How Effective is Hue’s Human Feedback Integration Mechanism?

For RQ3, we conducted experiments on all three hybrid logs as well as the three single-line logs that achieve relatively unsatisfactory grouping accuracy (below 0.9) and consist of a large number of log templates (over 100 templates) in Sec. 5.3 (i.e., Linux, Android, and BGL) to highlight the improvement brought by Hue’s feedback mechanism. Specifically, we set Hue to *guided updating* mode and evaluate the improvement on various datasets by changing the

maximum limit of triggered merge queries. Furthermore, in the experiments, we keep all feedback provided by the user correct to eliminate any potential perturbation caused by human factors.

The results are shown in Fig. 11 and Fig. 12. Aside from CTS, Hue’s grouping accuracy was able to improve to over 0.92 and template F1-score to over 0.87 for all datasets with less than 20 feedback queries. The performance on CTS might not be optimal, which could be due to the relatively high number of templates and low amount of corresponding logs in the dataset.

For effectiveness comparison, we selected SPINE-Feedback [32] to demonstrate the superiority of Hue’s feedback integration mechanism. We carried out experiments on Linux where both Hue and SPINE did not perform satisfactorily, i.e., with the lowest grouping accuracy of 0.749 and 0.676 among all 16 single-line datasets, respectively. To highlight the strength of Hue’s feedback mechanism, we intentionally reduced Hue’s base grouping accuracy by tuning its parameters to be lower than SPINE’s to highlight its feedback mechanism strength, despite its initially higher accuracy than SPINE’s in auto-updating mode.

The results are shown in Fig. 13. Compared with SPINE, Hue is able to improve accuracy efficiently with fewer feedback queries significantly. Moreover, it achieves a higher final grouping accuracy after tens of queries. The efficient and effective feedback performance of Hue may be attributed to two aspects: (1) Hue pre-filters a large number of unnecessary feedback queries by checking whether all tokens to be updated are keys. As discussed in Sec. 4.5.2, Hue assumes that keys are potential parameter tokens or non-ambiguous tokens. Therefore, the template updates that are solely targeting keys are considered as correct updates and will not trigger feedback queries. (2) The feedback mechanism of Hue resembles "*preventing erroneous template updates*," whereas the feedback mechanism of SPINE resembles "*correcting errors after they occur*." Specifically, SPINE iteratively partitions the already merged groups into new groups and triggers feedback queries each time. However, this may not be able to resolve the cases where logs belonging to the same template have already been clustered in the wrong groups, where the partition feedback queries cannot help them correctly clustered together again. In contrast, Hue raises feedback queries before merging logs into wrong groups, which can effectively prevent incorrect grouping cases at the very beginning. Thus, Hue can effectively achieve higher final accuracy.

Additionally, we observe that Hue’s accuracy finally saturates as the number of merge queries increases in some experiments. This is because the parsing errors could also come from other components of Hue (e.g., from the heuristic parse tree in fine-grouping).

## 6 INDUSTRIAL DEPLOYMENT

We deployed Hue in Huawei Cloud and it ran smoothly for two months, from October 2022 to December 2022. During this period, we annotated 1,000 logs for each downstream task every week. Throughout the evaluation, IT operators did not have access to the source code of the log statements.

To assess the performance of Hue, we calculated the percentage of each type of log message over the two-month period for root cause analysis tasks. As shown in Fig. 14, the blue part represents

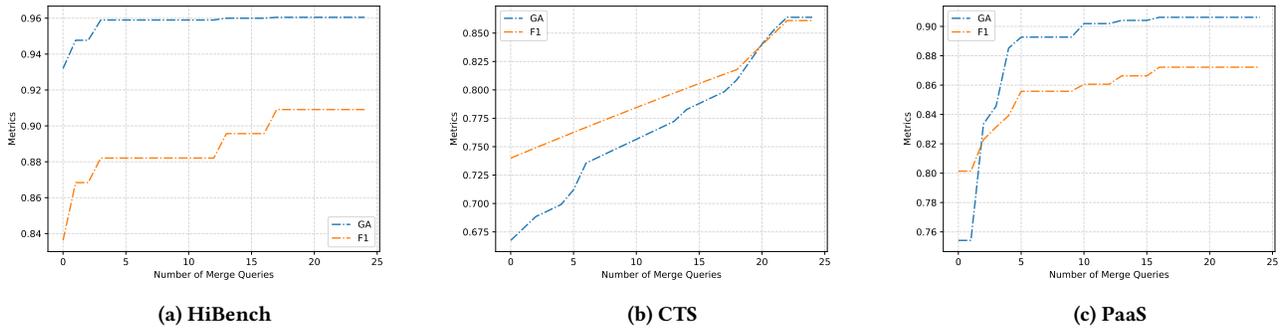


Figure 11: The increase of grouping accuracy and template F1-score on hybrid log datasets with feedback.

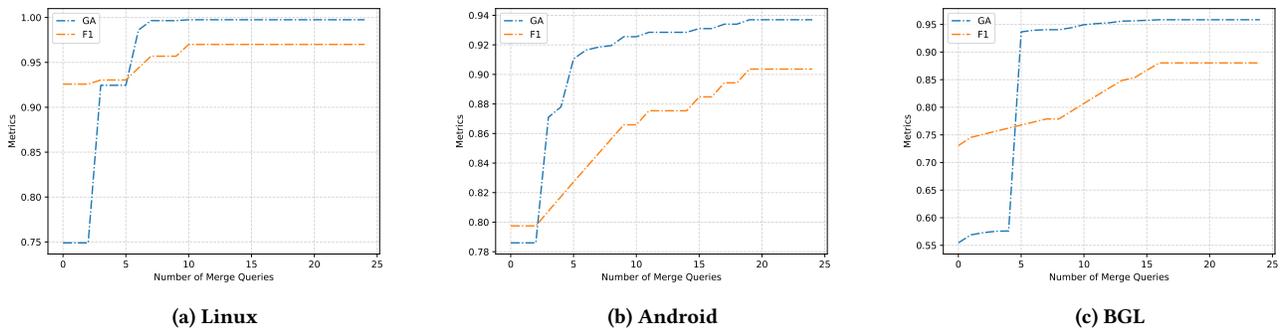


Figure 12: The increase of grouping accuracy and template F1-score on three LogPAI datasets with feedback.

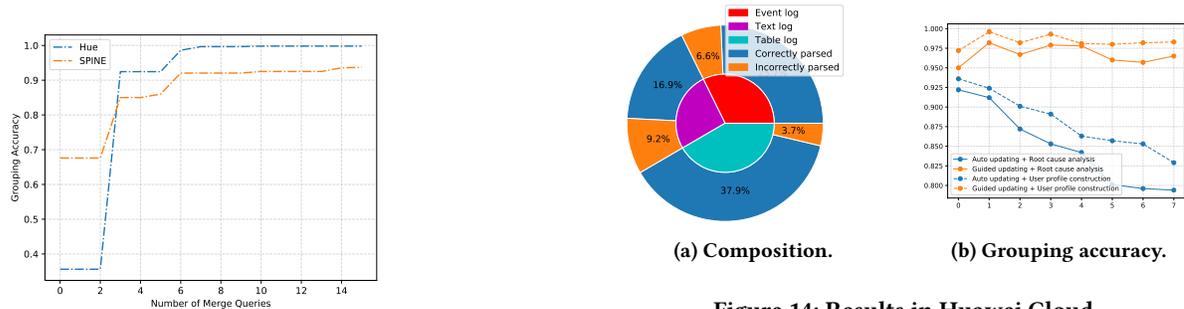


Figure 14: Results in Huawei Cloud.

Figure 13: Feedback mechanism comparison on Linux.

correctly parsed templates, while the orange part represents incorrectly parsed templates. Out of all parsed log templates, 32.3% were event logs, 26.1% were text logs, and 41.6% were table logs. Nearly 70% of them were multi-line, making it challenging for traditional log parsers to provide accurate log templates for downstream tasks.

In our department in Huawei Cloud, the logs are used for two downstream tasks: root cause analysis and user profile construction. Since IT operators cannot directly access the source code and extract templates from logging statements, they are limited to manually defining log templates according to their requirements. However, as discussed in Sec. 2.2, distinct downstream tasks may possess varying

interpretations of log templates. For example, event logs typically exhibit a "verb-object" structure, manifesting as "performed an action on a specific object" (Fig. 3). In the root cause analysis task, engineers focus on the objects being operated, so they partition templates based on objects. In the user profile construction task, engineers focus on the recorded action, so they partition templates based on actions. To achieve more accurate log parsing aligned with their downstream task requirements, they utilized Hue's guided updating mode to parse the collected log samples with different human feedback.

Fig. 14 shows the results, where the solid and dashed lines are for root cause analysis and user profile construction, and blue and

orange lines are for auto-updating and guided-updating modes, respectively. This performance has been deemed acceptable by the engineers responsible for their respective downstream tasks. They have successfully integrated the parsed hybrid log data into their experimental models. Furthermore, in our discussion with engineers, we learned that Hue did not trigger an excessive number of feedback queries during the actual parsing process, thus avoiding user fatigue in providing feedback. In fact, the average number of feedback queries triggered in every 1,000 logs was less than 30, while the grouping accuracy remained acceptable (*i.e.*, over 0.9). This enabled them to try to incorporate hybrid logs as available data in their models for the first time. We also found that Hue’s space-based tokenization strategy (discussed in Sec. 4.2) has minimal impact on the downstream tasks in Huawei Cloud, contrary to our previous belief that it could threaten the practicality of Hue. This is primarily due to the fact that Huawei Cloud’s business solely uses vectorized logs for various tasks, where the main requirement is accurate clustering of logs belonging to the same template, rather than ensuring every token in the extracted templates (represented as strings) is correctly parsed. Therefore, we believe that Hue also demonstrates a certain level of practicality in industrial scenarios.

## 7 THREATS TO VALIDITY

In this work, we identified the following major threats to validity.

- **Data Quality:** Our evaluation uses multiple log datasets and we found that they contain labeling errors that can negatively impact parser performance. Also, most parsers, including Hue, assume that logs have a specific header. Hue segments logs based on headers. For logs without headers, Hue’s performance would be degraded. Thus, we regard parsing hybrid logs without headers as future work.
- **Key Adaptability:** In Sec. 5.3, we use the most common key setting for all experiments. However, some logs may not fit this setting, which could affect the parsing result. To mitigate this threat, we allow users to customize their own key casting configurations.
- **Query Correctness:** In Sec. 5.4, we assume that human feedback to the query is always correct. However, if the user provided incorrect feedback, it might introduce new parsing errors. To mitigate this impact, we recommend that users discuss the query in detail if they are not sure about the correct template.

## 8 RELATED WORKS

Log parsing methods can be categorized into unsupervised and supervised methods according to the parsing algorithms.

**Unsupervised log parsers** do not require labels on existing log data and thus they have been more widely explored. There are three main categories of unsupervised log parsers: frequent pattern mining-based, clustering-based, and heuristic-based methods. (1) *Frequent pattern mining-based* methods regard the mined frequent patterns (*e.g.*, *n*-grams) as log templates. For example, SLCT [30], LFA [25], LogCluster [31], and Logram [2] try to use different methods to extract frequent patterns in logs. (2) *Clustering-based* methods aim to group similar logs first, assuming that logs in the same cluster share the same template, and then extract the

common tokens as the template in each cluster. Some clustering-based methods can perform in an online manner because they adopt an online grouping strategy rather than clustering all the offline logs at once. Specifically, LKE [8], LogSig [29], LogMine [10] are offline methods, SHISO [24], and LenMa [28] are online methods. (3) *Heuristic-based* methods encode expert domain knowledge into general and effective heuristic rules. For example, AEL [7], IPLoM [23], Spell [3], and Drain [12] utilize different heuristic rules to extract templates from logs. In particular, Drain [12] achieved SOTA in all open-source traditional parsers with a parse tree structure to perform log parsing in an online manner. POP [11] improves Drain and provides a parallel implementation on Spark for distributed deployment. SPINE [32] improves Drain and proposes a progressive clustering step for human feedback. In particular, SPINE introduced a manual correction paradigm to log parsing. Although SPINE also considers human feedback, this paper proposed a new human feedback integration mechanism based on the merge-reject strategy, which largely improves efficiency and achieved higher accuracy.

**Supervised log parsers** are less common because they rely on lots of precise ground truth labels manually constructed by IT operators. UniParser [21] utilizes a contrast learning strategy to overcome the pattern difference in heterogeneous log sources and uses a BiLSTM [9] model for token and context encoding. Sem-Parser [16] utilizes a special semantic miner for template extracting and uses a BiLSTM [9] model for encoding. Different from these supervised approaches, Hue is an unsupervised log parser that does not rely on label log datasets.

Different from these existing parsers, Hue is the first attempt to parse hybrid logs, providing an approach to utilize neglected hybrid logs for various downstream tasks in log analysis. Moreover, most existing parsers can only improve parsing performance through hyperparameter tuning, while Hue, in addition to hyperparameter tuning, further enhances parsing effectiveness through a feedback mechanism.

## 9 CONCLUSION

This paper introduces a new, general, and practical research problem, hybrid log parsing, which is the superset of the widely-studied single-line log parsing problem. We also present *Hue*, the first attempt to tackle this problem. Hue leverages both patterns in the incoming log messages and expert domain knowledge to parse both hybrid and single-line logs in an online manner. The use of a key casting table and a merge-reject strategy enables Hue to effectively utilize user feedback, making it possible to quickly adapt to complex and changing log templates. The results of our evaluation on three hybrid log datasets and sixteen widely-used single-line log datasets demonstrate the superiority of Hue in terms of accuracy and efficiency. The successful deployment of Hue in a real production environment further highlights its practical value. We hope Hue and its replication package can serve as the first step in hybrid log parsing and benefit a line of future research in this interesting and practice direction.

## 10 DATA AVAILABILITY

We uploaded our repository at <https://github.com/logpai/hybridlog>

## REFERENCES

- [1] Anunay Amar and Peter C Rigby. 2019. Mining historical test logs to predict bugs and localize faults in the test logs. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 140–151.
- [2] Hetong Dai, Heng Li, Che Shao Chen, Weiyi Shang, and Tse-Hsun Chen. 2020. Logram: Efficient log parsing using n-gram dictionaries. *IEEE Transactions on Software Engineering* (2020).
- [3] Min Du and Feifei Li. 2016. Spell: Streaming parsing of system event logs. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 859–864.
- [4] Min Du, Feifei Li, Guineng Zheng, and Vivek Srikumar. 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*. 1285–1298.
- [5] ELK. 2023. Manage multiline messages. <https://www.elastic.co/guide/en/beats/filebeat/current/multiline-examples.html#multiline-examples>. [Online; accessed 15-June-2023].
- [6] Fluentd. 2023. Multiline. <https://docs.fluentd.org/parser/multiline>. [Online; accessed 15-June-2023].
- [7] Qiang Fu, Jian-Guang Lou, Yi Wang, and Jiang Li. 2009. Execution anomaly detection in distributed systems through unstructured log analysis. In *2009 ninth IEEE international conference on data mining*. IEEE, 149–158.
- [8] Qiang Fu, Jian-Guang Lou, Yi Wang, and Jiang Li. 2009. Execution anomaly detection in distributed systems through unstructured log analysis. In *2009 ninth IEEE international conference on data mining*. IEEE, 149–158.
- [9] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 6645–6649.
- [10] Hossein Hamooni, Biplob Debnath, Jianwu Xu, Hui Zhang, Guofei Jiang, and Abdullah Mueen. 2016. Logmine: Fast pattern recognition for log analytics. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. 1573–1582.
- [11] Pinjia He, Jieming Zhu, Shilin He, Jian Li, and Michael R Lyu. 2017. Towards automated log parsing for large-scale log data analysis. *IEEE Transactions on Dependable and Secure Computing* 15, 6 (2017), 931–944.
- [12] 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)*. IEEE, 33–40.
- [13] Shilin He, Pinjia He, Zhuangbin Chen, Tianyi Yang, Yuxin Su, and Michael R Lyu. 2021. A survey on automated log analysis for reliability engineering. *ACM computing surveys (CSUR)* 54, 6 (2021), 1–37.
- [14] Shilin He, Qingwei Lin, Jian-Guang Lou, Hongyu Zhang, Michael R Lyu, and Dongmei Zhang. 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*. 60–70.
- [15] Shilin He, Jieming Zhu, Pinjia He, and Michael R Lyu. 2016. Experience report: System log analysis for anomaly detection. In *2016 IEEE 27th international symposium on software reliability engineering (ISSRE)*. IEEE, 207–218.
- [16] Yintong Huo, Yuxin Su, Baitong Li, and Michael R Lyu. 2021. SemParser: A Semantic Parser for Log Analysis. *arXiv preprint arXiv:2112.12636* (2021).
- [17] Zanis Ali Khan, Donghwan Shin, Domenico Bianculli, and Lionel Briand. 2022. Guidelines for Assessing the Accuracy of Log Message Template Identification Techniques. In *Proceedings of the 44th International Conference on Software Engineering (ICSE'22)*. ACM.
- [18] Kenneth Wai-Ting Leung and Dik Lun Lee. 2009. Deriving concept-based user profiles from search engine logs. *IEEE Transactions on knowledge and data engineering* 22, 7 (2009), 969–982.
- [19] Qingwei Lin, Hongyu Zhang, Jian-Guang Lou, Yu Zhang, and Xuewei Chen. 2016. Log clustering based problem identification for online service systems. In *2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C)*. IEEE, 102–111.
- [20] Jinyang Liu, Jieming Zhu, Shilin He, Pinjia He, Zibin Zheng, and Michael R Lyu. 2019. Logzip: extracting hidden structures via iterative clustering for log compression. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 863–873.
- [21] Yudong Liu, Xu Zhang, Shilin He, Hongyu Zhang, Liqun Li, Yu Kang, Yong Xu, Minghua Ma, Qingwei Lin, Yingnong Dang, et al. 2022. UniParser: A Unified Log Parser for Heterogeneous Log Data. In *Proceedings of the ACM Web Conference 2022*. 1893–1901.
- [22] Jian-Guang Lou, Qiang Fu, Shenqi Yang, Ye Xu, and Jiang Li. 2010. Mining invariants from console logs for system problem detection. In *2010 USENIX Annual Technical Conference (USENIX ATC 10)*.
- [23] Adetokunbo AO Makanju, A Nur Zincir-Heywood, and Evangelos E Milios. 2009. Clustering event logs using iterative partitioning. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1255–1264.
- [24] Masayoshi Mizutani. 2013. Incremental mining of system log format. In *2013 IEEE International Conference on Services Computing*. IEEE, 595–602.
- [25] Meiyappan Nagappan and Mladen A Vouk. 2010. Abstracting log lines to log event types for mining software system logs. In *2010 7th IEEE Working Conference on Mining Software Repositories (MSR 2010)*. IEEE, 114–117.
- [26] NETWITNESS. 2023. Multi Line Logs to Single Line Logs. <https://community.netwitness.com/t5/netwitness-discussions/multi-line-logs-to-single-line-logs-conversion/td-p/511346>. [Online; accessed 15-June-2023].
- [27] OBSERVE. 2023. Multi-line Logging Strategies. <https://www.observeinc.com/blog/multi-line-logging-strategies/>. [Online; accessed 15-June-2023].
- [28] Keichi Shima. 2016. Length matters: Clustering system log messages using length of words. *arXiv preprint arXiv:1611.03213* (2016).
- [29] Liang Tang, Tao Li, and Chang-Shing Perng. 2011. LogSig: Generating system events from raw textual logs. In *Proceedings of the 20th ACM international conference on Information and knowledge management*. 785–794.
- [30] Risto Vaarandi. 2003. A data clustering algorithm for mining patterns from event logs. In *Proceedings of the 3rd IEEE Workshop on IP Operations & Management (IPOM 2003)*(IEEE Cat. No. 03EX764). IEEE, 119–126.
- [31] Risto Vaarandi and Mauno Pihelgas. 2015. Logcluster—a data clustering and pattern mining algorithm for event logs. In *2015 11th International conference on network and service management (CNSM)*. IEEE, 1–7.
- [32] Xuheng Wang, Xu Zhang, Liqun Li, Shilin He, Hongyu Zhang, Yudong Liu, Lingling Zheng, Yu Kang, Qingwei Lin, Yingnong Dang, et al. 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*. 1198–1208.
- [33] Junyu Wei, Guangyan Zhang, Yang Wang, Zhiwei Liu, Zhanyang Zhu, Junchao Chen, Tingtao Sun, and Qi Zhou. 2021. On the Feasibility of Parser-based Log Compression in {Large-Scale} Cloud Systems. In *19th USENIX Conference on File and Storage Technologies (FAST 21)*. 249–262.
- [34] Shenglin Zhang, Ying Liu, Weibin Meng, Zhiling Luo, Jiahao Bu, Sen Yang, Peixian Liang, Dan Pei, Jun Xu, Yuzhi Zhang, et al. 2018. Prefix: Switch failure prediction in datacenter networks. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 2, 1 (2018), 1–29.
- [35] Xu Zhang, Yong Xu, Qingwei Lin, Bo Qiao, Hongyu Zhang, Yingnong Dang, Chunyu Xie, Xinsheng Yang, Qian Cheng, Ze Li, et al. 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*. 807–817.
- [36] Jieming Zhu, Shilin He, Jinyang Liu, Pinjia He, Qi Xie, Zibin Zheng, and Michael R Lyu. 2019. Tools and benchmarks for automated log parsing. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 121–130.