

Self-Evolutionary Group-wise Log Parsing Based on Large Language Model

Anonymous Author(s)

1 **Abstract**—Log parsing involves extracting appropriate templates from semi-structured logs, providing foundational information for downstream log analysis tasks such as anomaly detection and log comprehension. Initially, the task of log parsing was approached by domain experts who manually designed heuristic rules to extract templates. However, the effectiveness of these manual rules deteriorates when certain characteristics of a new log dataset do not conform to the pre-designed rules. To address these issues, introducing large language models (LLM) into log parsing has yielded promising results. Nevertheless, there are two limitations: one is the reliance on manually annotated templates within the prompt, and the other is the low efficiency of log processing. To address these challenges, we propose a self-evolving method called *SelfLog*, which, on one hand, uses similar `<group, template>` pairs extracted by LLM itself in the historical data to act as the prompt of a new log, allowing the model to learn in a self-evolution and labeling-free way. On the other hand, we propose an N-Gram-based grouper and log hitter. This approach not only improves the parsing performance of LLM by extracting the templates in a group-wise way instead of a log-wise way but also significantly reduces the unnecessary calling to LLMs for those logs whose group template is already extracted in history. We evaluate the performance and efficiency of *SelfLog* on 16 public datasets, involving tens of millions of logs, and the experiments demonstrate that *SelfLog* has achieved state-of-the-art (SOTA) levels in 0.975’s GA, and 0.942’s PA. More importantly, without sacrificing accuracy, the processing speed has reached a remarkable 45,000 logs per second.

29 **Keywords**—large language model, log parsing, self-evolution

30 I. INTRODUCTION

31 Logs [1] and time series [2], along with trace [3], jointly constitute the three vital types of data for monitoring and analyzing the reliability of software systems. Log data is the most easily acquired and the most widely encompassing. However, due to its semi-structured nature, it poses a greater challenge for analysis. Modern software systems such as operating systems like Windows and Android, or file systems like HDFS [1], generate a substantial volume of logs daily. Analyzing and detecting anomalies in such a vast number of logs manually is impractical. To achieve automatic log processing, **log parsing**, which involves transforming semi-structured logs into a structured format, making it the most crucial preliminary step for downstream tasks such as log anomaly detection [4], [5], log compression [6], [7], and log summarizing [8].

46 As shown in Fig. 1, log parsing primarily involves distinguishing between the constant and variable parts (named parameters in Fig. 1) of the log content through a series of methods. By replacing the variable parts with wildcards, we create a log template. It can be observed from Fig. 1 that

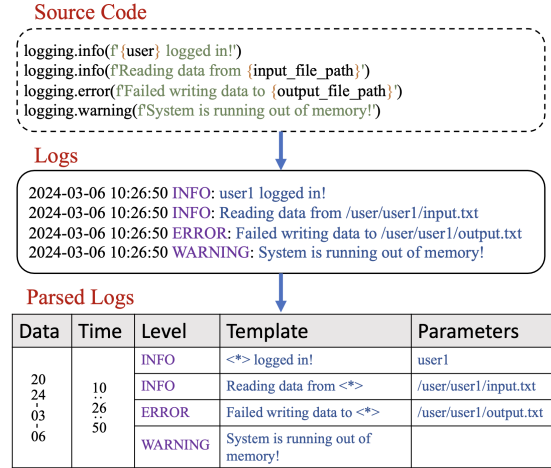


Fig. 1. An illustration example of log parsing.

51 if the source code is available, it becomes much easier to extract the corresponding template. However, in many systems, 52 the original software code is not accessible, leading to the 53 emergence of many data-driven log parsing methods [9]– 54 [11]. These data-driven methods are primarily categorized 55 into unsupervised and supervised approaches. Unsupervised 56 methods typically employ heuristic rules [12]–[14] or sta- 57 tistical features [9]–[11] to extract templates. However, the 58 effectiveness of these methods can be greatly impacted if the 59 log datasets to be analyzed do not align well with the pre- 60 designed rules or features. For example, in Drain [12], it is 61 assumed that the first token of each log is constant. Yet, in 62 real-world systems, we have found that this is not always the 63 case, such as with “proxy.cse.cuhk.edu.hk:5070 64 close, 451 bytes sent, 353 bytes received, 65 lifetime <1 sec”, where the first token is variable. 66 Supervised log parsing methods [15], [16], on the other hand, 67 train or fine-tune models using manually annotated <log, 68 template> pairs or token types. However, these methods are 69 sensitive on the distribution of the training data and may 70 perform poorly on logs unseen during training. 71

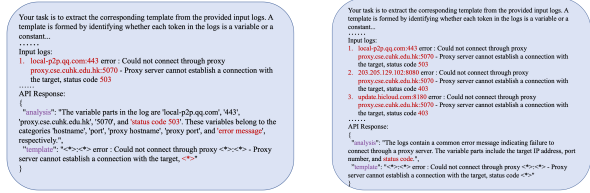
72 To address the issues with the aforementioned methods, 73 some current research has begun to employ Large Language 74 Models (LLMs) for log parsing [17]. Logs are inherently 75 statements printed by programmers, naturally containing se- 76 mantic information, and LLMs are inherently capable of 77 extracting semantic information. Furthermore, LLMs, due to 78 their powerful zero-shot capabilities, can better transfer to a

79 new set of logs without the need for additional hyperparameter
 80 adjustments. Thus, using LLMs for log parsing is a promising
 81 direction. However, existing LLM-based log parsing methods
 82 have the following two drawbacks:

- 83 • Current methods, such as DivLog [17], provide the LLM
 84 with prompts containing similar logs and corresponding
 85 templates, enabling the LLM to extract new log tem-
 86 plates through in-context learning (ICL). However, this
 87 approach heavily relies on the quality of the examples
 88 in the prompt. When software systems undergo upgrades
 89 and iterations, a new round of manual annotation of new
 90 log templates is required.
- 91 • Most importantly, existing LLM-based methods do not
 92 explicitly evaluate and discuss the log processing speed
 93 and token cost after deployment, which is key to de-
 94 termining whether LLMs can be widely applied in log
 95 parsing. We find through actual measurement that existing
 96 LLM-based methods can only process no more than 10
 97 logs per second (as detailed in Section IV-F), whereas
 98 in real software systems, it is very common for tens of
 99 thousands of logs to be generated per second.

100 To address these challenges, we propose a self-evolutionary
 101 group-wise LLM-based log parsing system called *SelfLog*.
 102 Although the templates extracted by the LLM from historical
 103 logs are not as accurate as those of domain experts, they can
 104 provide useful information for LLM to revise its responses.
 105 Inspired by this, we build a *Prompt Database* to store LLM
 106 history extracted <log, template> pairs. When a new log
 107 arrives, the most similar logs can be retrieved through an
 108 approximate nearest neighbor (ANN) search. Then, by incor-
 109 porating similar logs and the templates previously extracted
 110 by the LLM itself into the prompt, it can help the LLM
 111 reflect on the correctness of the previously extracted templates,
 112 thereby improving the extraction performance for new logs
 113 and helping *SelfLog* no longer rely on manual annotation by
 114 human experts.

115 In the meantime, we observe that existing LLM methods
 116 process logs one at a time for template extraction, as depicted
 117 in Fig. 2(a), a method we refer to as point-wise parsing.
 118 However, this approach overlooks certain information. Typ-
 119 ically, domain experts determine the variables in a log by
 120 comparing it with several similar logs and identifying the
 121 differences between them. As shown in Fig. 2(a), if a single
 122 log is given to an LLM, it might incorrectly identify `status`
 123 `code 503` as a variable. Yet, if we present a group of
 124 similar logs to the LLM, as shown in Fig. 2(b), the model
 125 accurately identifies 503 as the variable and `status code`
 126 as a constant, resulting in the correct template: `status`
 127 `code <*>`. Therefore, we design an *N-Gram-based Grouper*
 128 that first groups the logs and then invokes the LLM to extract
 129 templates for each group. This method not only enhances
 130 accuracy but also has the advantage of reducing the need to
 131 invoke the LLM for each log. We only need to call the LLM to
 132 extract templates for each group. For groups whose templates
 133 have already been extracted, we store the templates in our



(a) Point-wise parsing. (b) Group-wise parsing.

Fig. 2. An example from the Proxifier dataset [1], demonstrating the template extraction effectiveness of LLM when each log is given to the model, versus when several similar logs are grouped together and then given to the model.

134 designed *Log Hitter*. Groups that hit *Log Hitter* return directly
 135 without invoking the LLM, which significantly increases the
 136 parsing speed of *SelfLog* and substantially reduces the number
 137 of tokens required for model calls.

138 We evaluated *SelfLog* on 16 public datasets, and its perfor-
 139 mance exceeded the current state-of-the-art (SOTA). In Group
 140 Accuracy, its precision was 10% higher than SOTA, and in
 141 Parsing Accuracy, it improved by 16%. At the same time,
 142 we tested it on a dataset of tens of millions of logs, and our
 143 processing rate reached 45,000 logs per second, meeting the
 144 log generation rate of current online systems.

145 In summary, our main contributions are as follows:

- 146 • We propose *SelfLog*, a self-evolutionary group-wise log
 147 parser that achieves excellent log parsing results without
 148 the need for manual annotation of new templates, relying
 149 solely on LLM’s historical parsing results through con-
 150 tinuous reflection and correction.
- 151 • For the first time, we focus on the efficiency issue
 152 in LLM-based log parsing. By combining a specially
 153 designed N-Gram-based Grouper, Log Hitter, *SelfLog* can
 154 parse over 45,000 logs per second with higher parsing
 155 accuracy. At the same time, *SelfLog* only consumes 1%
 156 tokens quota compared with the SOTA LLM-based log
 157 parsing method, making it affordable for real-world sys-
 158 tems.

159 The following sections of the paper are organized as fol-
 160 lows: In Section II, we introduce the motivation. In Section III,
 161 we detail the implementation of our *SelfLog*. In Section IV,
 162 we describe the experimental setup and evaluate the algorithm.
 163 Section V discusses threats to validity. Section VI reviews
 164 related work, and Section VII summarizes the paper.

165 II. MOTIVATION AND BACKGROUND

166 In this section, we present a broad definition of system logs
 167 and explain the basic steps involved in log parsing. We also
 168 introduce the large language models (LLM), as well as the
 169 background knowledge related to in-context learning (ICL)
 170 and prompts.

171 A. Log Parsing

172 The system generates a large amount of logs every day [18],
 173 [19]. It is unrealistic to rely on operation and maintenance
 174 person to manually check the logs to detect system abnor-
 175 malities in time. Therefore, fully automatic log processing

176 is necessary [1], [20]. The first step in log processing is
177 log parsing, which processes semi-structured and unstruc-
178 tured logs and converts them into structured data for other
179 downstream tasks such as log anomaly detection [4], [5], log
180 compression [6], [7], and root cause analysis [21], [22]. Log
181 parsing includes preprocessing and log template extraction,
182 as shown in Fig. 1. The original log includes the timestamp
183 automatically stamped by the system, verbosity level, and log
184 content written in the program code. The formats of the first
185 two parts are bound to the system settings and only require
186 fixed regular expressions to be aligned and extracted. The
187 log content consists of a constant string written in the code
188 and a part that changes dynamically according to the system
189 status. The log parser extracts the constant parts from the log
190 content and replaces the variable parts with wildcards, such
191 as the asterisk (*). This text string, composed of constants
192 and wildcards, is commonly referred to as a log template.
193 A single log template corresponds to multiple logs where
194 the variables take on different values. The earliest log parser
195 directly parsed the system code and extracted the log template
196 according to the log output statement [23], [24] which is
197 shown in the source code part of Fig. 1. In many scenarios,
198 the original source code of the system itself is not accessible,
199 hence there is a substantial amount of work that relies on
200 extracting templates directly from the output logs themselves,
201 which can be divided into unsupervised method [9], [10], [12]
202 and supervised method [15], [16], [25]. However, there are
203 still various shortcomings and the effectiveness need to be
204 improved.

205 Unsupervised log parsers, which propose heuristic rules to
206 extract log templates based on the author’s observation of logs,
207 often have design flaws that cannot correctly parse all logs.
208 For example, when Drian [12] clusters logs, it first divides
209 them according to the length after word segmentation. How-
210 ever, in the publicly available dataset Proxifier [1], there are
211 two log entries: `...received, lifetime <1 sec` and
212 `...received, lifetime 00:02`. After tokenization by
213 whitespace, these two logs will have different lengths, but they
214 actually belong to the same template. Furthermore, Drain [12]
215 assumes that the first token of each log is a constant, but in
216 Proxifier, there are logs that start with variables. These exam-
217 ples illustrate that existing methods based on heuristic rules
218 or statistical features require special treatment for different
219 datasets. If the dataset is not properly handled, the accuracy
220 of template extraction can be greatly reduced.

221 Supervised log parsing requires training a template extrac-
222 tion model, or fine-tuning a pre-trained model, based on a
223 dataset that has been manually annotated. However, manual
224 annotation is costly, especially when the system generating
225 the logs is dynamically changing and undergoing updates
226 and upgrades, making manual annotation even more difficult.
227 VALB [25] trains a BiLSTM [26] model by manually annotat-
228 ing variable types to distinguish between the types of constants
229 and variables. They have predefined nine types of variables,
230 such as Object ID (OID) and Location Indicator (LOI). In
231 real-world datasets, such as LogPai [1], some corresponding

232 templates have very few logs associated with them, which is
233 insufficient for model training or fine-tuning.

B. Large Language Model and In-Context Learning

234 Large Language Models (LLMs) are a subcategory of
235 machine learning models. Initially, they are used in the field
236 of Natural Language Processing (NLP) [27] to understand and
237 generate text that is readable by humans. Subsequently, their
238 application has been extended to include images (PLEASE
239 add reference) and speech (PLEASE add reference). The
240 widespread use of LLMs is not only due to their large number
241 of parameters but, most importantly, their impressive ability to
242 follow instructions, which allows them to be employed for a
243 variety of different downstream tasks and demonstrates their
244 strong zero-shot capabilities. These tasks include translation,
245 text generation, and logical reasoning, among others.

246 In-context learning (ICL) is a vital aspect of LLMs [28].
247 In the context of Large Language Models (LLMs) such as
248 GPT-3 [29], a prompt represents the initial user input that
249 kickstarts the process of text generation by the model. The
250 nature and specificity of the prompt broadly determine the
251 direction and content of the generated response [30]. Prompts
252 can range from single words to complex paragraphs. Through
253 their training on diverse and extensive data, LLMs can handle
254 a wide spectrum of prompts, generate appropriate responses,
255 and provide insights or narratives based on them. In essence,
256 the prompt is the steering wheel that guides the text generation
257 journey of the LLM. The ability to absorb, interpret, and
258 utilize the shared contextual information in the conversation
259 to provide more relevant and useful responses [31] of LLM is
260 called ICL.

261 The emergence of Large Language Models (LLMs) and
262 Inductive Conformal Logic (ICL) has provided new insights
263 for robust and universal log parsing. System logs are output
264 to log files by programmers through code. They are records of
265 events that happen within a software application, system, or
266 network [24]. To facilitate operation and maintenance people
267 to monitor the system based on the logs, the logs are highly
268 readable and very similar to natural language. Besides, there
269 is also log information in the training data of LLM [32]. With
270 the human knowledge and in-context examples in prompts,
271 DivLog [17] found that it can achieve better results than SOTA
272 log parsers in log template extraction without fine-tuning
273 LLM. In real scenarios, the number and content of system
274 log templates change dynamically. Traditional log parsers may
275 need to re-train or re-design rules to improve the parsing effect
276 of new logs. With the powerful natural language understanding
277 capabilities of LLM, LLM-based Log Parsers only need to
278 change the in-context examples in prompt to achieve accurate
279 parsing of new log.

280 However, even though LLMs can effectively understand
281 logs, their calling costs and the speed at which they generate
282 templates are quite limited. In industrial systems, it is very
283 common to produce tens of thousands of logs in a single
284 minute, which presents a challenge for large models to process
285 within such a short time frame. In Section III, we design
286

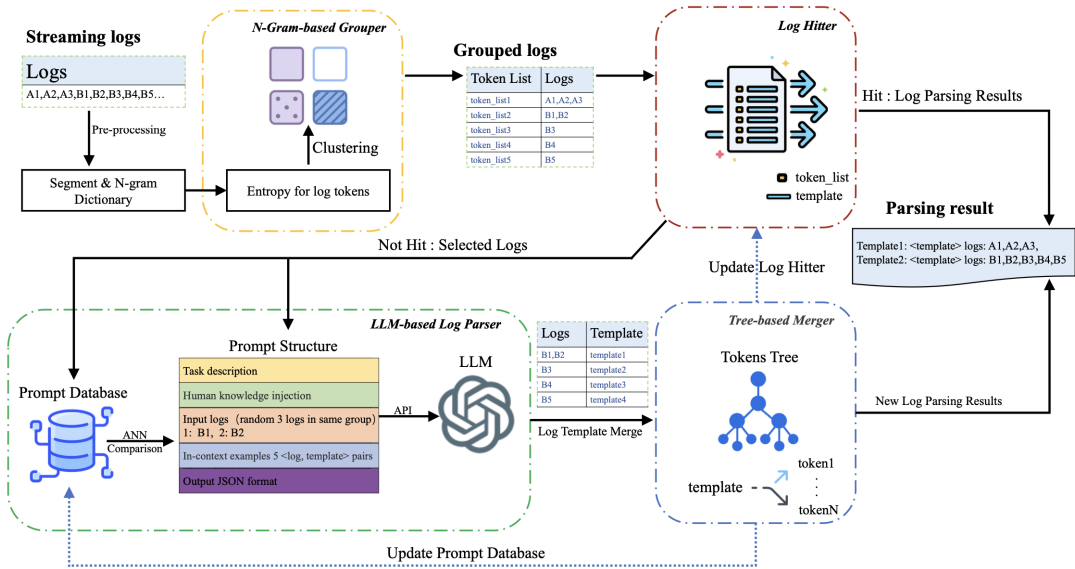


Fig. 3. The workflow of *SelfLog* framework. Here we demonstrate how to process streaming logs, but this system is also capable of seamlessly handling offline logs. In the case of offline log processing, only a forward pass (represented by the solid black line in the diagram) is necessary, without the need for iterative updates to the online database operations (represented by the dashed blue lines in the diagram).

287 algorithm that utilizes a self-evolution approach to generate
 288 higher-quality prompts without the need for manual expert
 289 annotation, thus enhancing the effectiveness of log parsing.
 290 Concurrently, we have also designed a framework to increase
 291 the efficiency of log processing using LLMs, aiming to in-
 292 crease the log processing rate of existing LLM-based log
 293 parsing systems(ADD CITATION) from one log per second
 294 to tens of thousands of logs per second.

295 III. OUR APPROACH

296 In this section, we present *SelfLog*, a self-evolutionary
 297 log parsing system that focuses on how, through a self-
 298 evolutionary approach, Large Language Models (LLMs) can
 299 achieve good results without relying on manually annotated
 300 <log, template> data. Additionally, it addresses the issue
 301 where the log processing efficiency is influenced by the gener-
 302 ation rate of the LLM. We first provide an overview of *SelfLog*,
 303 followed by a detailed introduction to each of its key modules.
 304 It is noteworthy that we will introduce the system with a
 305 focus on the most common scenario in industrial systems:
 306 streaming logs. In a streaming setup, logs are continuously
 307 generated, and downstream log parsing algorithms need to
 308 extract templates in real-time. At the end of this section, we
 309 will discuss how the scenario of offline analysis is actually a
 310 special case of online streaming analysis.

311 A. Overview of *SelfLog*

312 As shown in Fig. 3, the *SelfLog* system only needs to
 313 call the LLM API, through In-Context Learning (ICL), by
 314 providing a specially designed prompt to the LLM, to perform
 315 the task of log parsing without the need to train the LLM. The
 316 entire system is divided into four major modules: **N-Gram-
 317 based Grouper, Log Hitter, LLM-based Log Parser, and
 318 Tree-based Merger.**

The N-Gram-based Grouper is primarily used to cluster and
 group the preprocessed logs, which has two main goals. One
 is that it can greatly save on the financial cost of calling the
 large model, as we no longer need to call the model for each
 log entry. Instead, we only need to call the model once for a
 new group as a whole. Additionally, it can greatly enhance
 accuracy because if there is only one log, the model can
 only determine which token is a variable based on semantic
 information. However, if there are multiple logs within the
 same group, the model can determine which token is a variable
 by comparing which parts of the logs differ, which is also the
 core idea behind methods like Drain [12]. We are the first to
 apply this to LLMs.

With the Grouper in place, the design of the Log Hitter
 naturally follows, as the same group is likely to correspond to
 a same template. If this group already has an existing template,
 there is no need to make further calls to the large model,
 which is particularly important in an online environment. This
 is because repeated calls to the model not only greatly reduce
 processing efficiency but also, due to the hallucinations and
 uncertain outputs of large models, can lead to unstable parsing
 results. The LLM-based Log Parser mainly achieves good
 detection results through a carefully designed prompt and the
 method of ICL. The in-context examples in the prompt play an
 important role for the model to continuously revise its current
 output based on the previous outputs in a self-evolution way.
 The Tree-based Merger is primarily responsible for correcting
 the model's output because neither the Grouper nor the LLM
 can guarantee a 100% accuracy rate. The Merger will merge
 some logs that were incorrectly divided into multiple groups
 and templates, thereby enhancing the model's precision.

The original logs contain the timestamps assigned by the system to the log content, and log types such as INFO and ERROR, process ID, etc. In the same system, these contents are all in fixed locations in the log, so they can be extracted through simple regular expressions. For example, “17/06/09 20:10:40 INFO spark.SecurityManager: Changing view acls to: yarn,curi” is a raw log from Spark [1] system. The regular expression “[r'(\d+\.){3} \d+', r'\b[KGTM]?B\b', r'([\w-]+\.){2},[\w-]+'” can extract each part separately. Thus we only need to focus on the log content, which is printed by the system code through the log print statement (see Fig. 1), which contains fixed constants prewritten in the code and variables dynamically filled in based on system operating information.

Since methods based on statistical features [12] require calculating the characteristics of a segment of text within a log to determine whether that segment is a variable or a constant, how to segment a log text becomes critically important. Existing methods mostly involve segmenting a log by using delimiters, and converting it into separate segments of text, each of which we refer to as a token. However, this method of segmentation necessitates selecting appropriate delimiters for different datasets. For instance, for BGL [1] logs, the delimiters might be “. ()”, while for Windows datasets [1], the delimiters could be “=: []”. We believe that this approach to segmentation is not robust enough. With the advent of large models, we no longer need to strictly rely on such tokenization rules. In *SelfLog*, the purpose of tokenization during preprocessing is to facilitate subsequent log grouping, rather than directly extracting templates. Therefore, we only need to identify the commonalities across multiple logs and filter out as many of the variable parts as possible. Hence, we have chosen “[A-Za-z0-9*]+” as our tokenization rule. It can be seen that anything not composed of letters and numbers is used as a delimiter. For example, for the log “pam_unix(sshd:auth): authentication failure; logname= uid=0 euid=0 tty=ssh ruser= rhost=202.100.179.208 user=root”, after our tokenization process, it becomes a list of constant tokens “unix sshd auth authentication failure logname euid ruser rhost user root”. By comparing the token list with the original log, we can observe that pure numeric sequences have been removed, as we assume that pure numbers are most likely variables and cannot represent a category of logs. Furthermore, to mitigate the potential impact of variables, we query the WordNet [33] lexicon for all tokens that, after tokenization, are three characters or fewer in length. If a token appears infrequently in WordNet, we consider it to be an invalid word and likely a prefix or suffix. In the previous example, tokens such as “pam”, “uid”, “tty”, and “ssh” were eliminated.

Algorithm 1: N-Gram-based Grouper

Input : log X

- 1 $TX = \text{get_token_list}(X)$
- 2 // *step1: find 2-gram const token*
- 3 $\text{position} = \text{get_2gram_const_index}(TX)$
- 4 // *step2: Get variable token list from right part*
- 5 $\text{variable_list_right} = \text{PILAR_gram}(TX, \text{position})$
- 6 // *step3: Get variable token list from left part*
- 7 $\text{variable_list_left} = \text{PILAR_gram}(TX, \text{position})$
- 8 // *step4: removing variable from TX*
- 9 $CX = TX - \text{variable_list_right} - \text{variable_list_left}$

Output: CX : constant tokens of log X

C. N-Gram-based Grouper

In the pre-processing phase, each log is represented with a list of tokens, and yet some of these tokens may be variables. We need to further identify these variables, remove them from the token list, and then use the token list for grouping, ensuring that logs within each group belong to the same template. It is worth noting that even if we do not identify all variables here, leading to logs of the same template being divided into two different groups, it is not a problem because later stages involving the Large Language Model (LLM) and the Tree-based Merger can correct them. We have improved upon the entropy-based method from PILAR [10] to determine whether a token is a variable or a constant, with the specific method detailed in Algorithm 1. The constants in the log are written in the code by programmers to facilitate the person to observe the system status and code debugging. Therefore, constant tokens are often words with higher frequency in the corpus. Different tokens are assigned different frequencies according to their frequency of occurrence in WordNet [33]. Function `get_2gram_const_index` (Line 3) calculates the largest sum of 2 consecutive token weights and returns their position. Then, starting from the position and moving to the right (Line 5), the algorithm employs the method from PILAR, using a 3-gram approach to dynamically determine whether each token is a variable. Each token is based on the ratio of the number of co-occurrences with its neighbors and the number of neighbor occurrences after removing itself, compared with the set threshold. If it is less than the threshold, it is considered a variable. Function `PILAR_gram` is the algorithm in listing 1 of the PILAR. It returns the `variable_list_righ`. Similarly, starting from the position and moving to the left, it determines whether each token is a variable (Line 7). Finally, return CX (Line 9).

Compared to PILAR, our N-Gram-based Grouper differs in the following two ways: Firstly, PILAR relies on assuming that the first word of the log is a constant, but this is often inaccurate because some logs start with variables. By checking the log templates in the Proxifier ground truth, we found that 2000 logs all start with variables. Because the algorithm in

442 PILAR defaults to the first token as a constant, the execution
 443 direction of the algorithm is from left to right. We judge
 444 whether a token is a variable based on the weight value, the
 445 starting point of the algorithm needs to be executed in both
 446 directions from right to left and from left to right to calculate
 447 the entropy of log tokens. Secondly, Unlike PILAR, which
 448 sets thresholds based on expert experience, we directly set the
 449 threshold to be automatically adjusted according to the number
 450 of different logs to improve the robustness of the group stage.
 451 After obtaining the list of constant tokens for each log
 452 through Algorithm 1, we then categorize the logs into different
 453 groups based on the token list. Each group is keyed by the
 454 token list, with the value being a list of logs that records all
 455 logs belonging to that group. Subsequently, the LLM-based
 456 Log Parser will extract the corresponding log template for each
 457 newly emerged group, and after the template is extracted, it
 458 will be updated into the log hitter in the form of a <token_list,
 459 template> pair.

460 D. Log Hitter

461 After grouping, the logs are divided into multiple groups
 462 according to the token list. The Log Hitter maintains a
 463 dictionary with the token list as the key and the log template
 464 as the value. The grouped logs will first be looked up in
 465 the dictionary according to the token list, and if there is a
 466 hit, the corresponding template will be directly returned to
 467 complete the log parsing. If there is no hit, the token list will be
 468 recorded as the key first, and the three logs with large editing
 469 distances in the group will be selected as the input of the
 470 LLM-based Log Parser, and the logs will be parsed by LLM.
 471 Finally, the log template obtained after Tree-based Merger
 472 processing is updated to the dictionary. Log Hitter records
 473 historical grouping information and continuously updates it.
 474 Only logs that have not appeared before are handed over to
 475 LLM for processing, which greatly improves the efficiency of
 476 log parsing.

477 E. LLM-based Log Parser

478 A model prompt is a brief text snippet provided to an
 479 LLM model to guide its generation of related content. Unless
 480 stated otherwise, we use GPT-3.5 as our LLM model, and
 481 we also evaluate the performance of other LLMs in the
 482 evaluation section. These prompts are typically crafted as
 483 questions, descriptions, or instructions to elicit the model’s
 484 output on specific topics or styles. By cleverly constructing
 485 prompts, it’s possible to steer the model towards generating
 486 text that aligns with expectations, thereby meeting user needs
 487 or accomplishing particular tasks. In this paper, we carefully
 488 design prompts to guide LLM in log template extraction. As
 489 shown by the different colors in Fig. 4, our prompts mainly
 490 consist of the following five parts, which we will introduce
 491 one by one.

492 **Task Description:** This part should be placed at the very
 493 beginning of the prompt to clearly state the task that the LLM
 494 needs to perform, and it is part of the instruction section. Our
 495 specific task description is shown in the figure. In addition

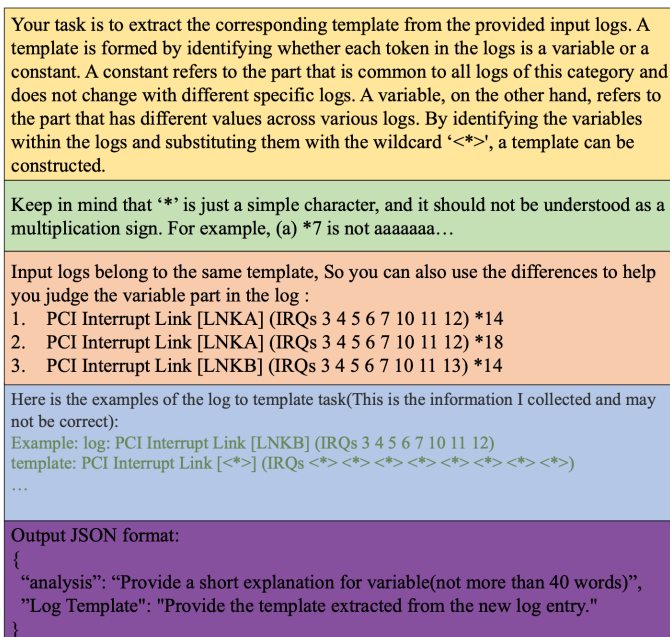


Fig. 4. An example of the complete prompt. The yellow block is the task description and the green block is human knowledge. The apricot block is the selected three input logs from the same group. The blue block is an example dynamically selected based on the Approximate Nearest Neighboring (ANN) from the prompt database, which is one of the core designs of our work. The purple part specifies the output format.

to providing the model with instructions for log extraction, we also inform the model of the system to which these logs belong, activating the corresponding log training part within the model.

Human Knowledge Injection: This part is optional. If there is explicit knowledge that can be articulated in actual applications, it can be added here to enhance the model’s expression. We include knowledge to inform the model that the asterisk (*) is not a multiplication sign but a representative of a wildcard, to prevent conflicts with other parts of the model’s knowledge. Examples of log machine correspondence templates based on historical manual confirmation in DivLog can also be placed in this part. Therefore, our algorithm can be combined with DivLog to achieve better improvement.

Input Logs: This part is the main input corresponding to the log template extraction task. Through the design of the N-Gram-based Grouper mentioned earlier, our model no longer extracts templates from individual logs. Instead, we extract new templates for each group. So, when a new group that has not been seen before appears, we randomly select three logs with the greatest edit distance from this new group as the model’s input for template extraction. In Fig. 2, we demonstrate the difference in final effect between using the LLM to perform template extraction on each log entry individually and feeding multiple similar logs into the LLM as a group for log parsing. From Fig. 2 (a), it can be seen that when a single log is fed to the LLM, the model will identify the status code 503 as a variable type of error message, thereby

524 recognizing the entire status code 503 as a variable. In
 525 fact, the status code is a constant, and 503 is the variable.
 526 If similar logs are input into the model as a group, as shown
 527 in Fig. 2(a), the model sees both the status code 503
 528 and the status code 403, thus accurately identifying the
 529 variable part. We will evaluate in Section 5 the impact of
 530 choosing different numbers of logs in a group on the final
 531 outcome.

532 **Self-evolution Examples:** This part is the main part de-
 533 signed in this paper. DivLog [17] works by adding manu-
 534 ally annotated logs and their corresponding templates to the
 535 prompt, but this still needs manual annotation for new logs.
 536 In this paper, we record the logs and their corresponding
 537 templates that the LLM has parsed in history, storing them in
 538 the Prompt Database. Each time a new log needs to be parsed,
 539 we retrieve the most similar historical logs and their templates
 540 from the data through an Approximate Nearest Neighbors
 541 (ANN) search, serving as the corpus for In-Context Learning
 542 (ICL). This approach not only allows for complete automation
 543 without the need for expert annotation but also enables the
 544 model to reflect on potential issues in previously extracted
 545 templates and make timely corrections.

546 **Output JSON Format:** The content provided by the LLM
 547 model is usually quite diverse and often includes some analysis
 548 and explanations of the problem. These outputs are usually
 549 mixed with the extracted templates. If there are no constraints
 550 on the model’s output, it would be difficult to directly extract
 551 the answer from the large model’s response. Therefore, we
 552 impose explicit constraints to let the LLM fill the analysis
 553 process and the final template into the pre-set json fields,
 554 which facilitates the subsequent accurate extraction of the log
 555 template from LLM’s answers.

556 F. Tree-based Merger

557 We test existing large models and find that even models
 558 like GPT-3 cannot identify all variables. As shown in the
 559 example in Fig. 5, since logs are mostly entered in a streaming
 560 manner, and the initial logs are all from the user “cyrus”, the
 561 model will extract a template with `session opened for`
 562 `user cyrus by (uid=<*>)`. Following this, when logs
 563 from the user “news” are entered, the model will propose a
 564 corresponding template, and so on. When there is a period
 565 with logs from both “cyrus”, “news”, “test”, and other users,
 566 the model can recognize that what lies between “user” and
 567 “by” is a variable. To address this issue, we construct a tree
 568 as depicted in Fig. 5. This tree updates in real-time based
 569 on the parsing results of the streaming data. By utilizing this
 570 tree, we can perform a double check of corner cases that the
 571 large model cannot accurately recognize, thereby enhancing
 572 the parsing performance of the model.

573 IV. EVALUATION

574 In this section, we design detailed experiments to answer
 575 and verify the following six research questions:

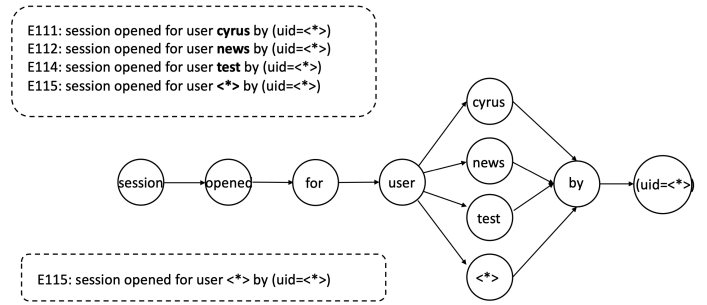


Fig. 5. Illustration of Tree-based Merger.

RQ1: Effectiveness of SelfLog. How does SelfLog perform 576
 in comparison to other state-of-the-art algorithms across the 577
 16 publicly annotated datasets by LogPai [1]? 578

RQ2: Efficiency and Cost of SelfLog. Compared with the 579
 LLM-based log parsing method, how efficient is SelfLog? 580

RQ3: Ablation Study. How do the different constituents in 581
 our design contribute to overall performance? 582

RQ4: Parameter Sensitivity. How do configuration parame- 583
 ters affect the parsing effects? 584

RQ5: Parsing Speed. What is the maximum parsing speed at 585
 which SelfLog currently processes streaming logs? 586

RQ6: LLM Backbone. What is the impact of different LLMs 587
 on the SelfLog effect? 588

A. Experimental setup 589

1) *Datasets:* The experimental dataset comes from the real 590
 log data of 16 different systems open-sourced by LogPai [1]. 591
 LogPai manually labeled templates of 2K logs for each dataset. 592
 593

2) *Evaluation metric:* Consistent with recent research find- 594
 ings [34], we employ Parsing Accuracy (PA), Precision Tem- 595
 plate Accuracy (PTA), and Recall Template Accuracy (RTA). 596
 Additionally, we have incorporated the Group Accuracy (GA) 597
 metric as used in the paper [12], [15], [16], [25], [34]. 598

- GA (Group Accuracy) was initially introduced by the pa- 599
 per [34] and has since been adopted for strictly assessing 600
 the accuracy of log template extraction. It considers a 601
 template extraction to be correct only if all correspond- 602
 ing logs belonging to the same template are accurately 603
 extracted. 604
- PA (Parsing Accuracy) was first proposed by LogGram 605
 [9]. PA focuses on the consistency between the log 606
 template extracted by the algorithm and the ground truth. 607
 If all tokens in the log are correctly identified as constants 608
 and variables, the extraction is considered correct. 609
- PTA (Precision Template Accuracy) and RTA (Recall 610
 Template Accuracy) is proposed by Khan et al. [34]. 611
 PTA is measured by the percentage of correctly identified 612
 templates to the total number of identified templates, 613
 whereas RTA is measured by the percentage of correctly 614
 identified templates to the total number of ground truth 615
 templates. 616

TABLE I

ACCURACY COMPARISON WITH DIFFERENT LOG PARSERS ON LOGPAI DATASETS ([1]). THE BEST SCORES FOR EACH METRIC OF EVERY DATASET ARE BOLD. DUE TO LIMITED TABLE SPACE, WE OMIT PTA AND RTA BECAUSE THEY SHOW CONSISTENT RESULTS WITH PA. IT IS NOTEWORTHY THAT, IN ADDITION TO THIS TABLE, WE INCLUDE GA, PA, PTA, AND RTA IN THE FOLLOWING FIGURES AND TABLES.

Dataset	LenMa		Spell		Drain		Logram		LogPPT		DivLog		SelfLog	
	GA	PA	GA	PA	GA	PA	GA	PA	GA	PA	GA	PA	GA	PA
HDFS	0.998	0.01	1.000	0.297	0.998	0.3545	0.940	0.005	0.845	0.389	0.143	0.966	1.000	1.000
BGL	0.690	0.082	0.787	0.197	0.963	0.342	0.645	0.125	0.478	0.789	0.451	0.949	0.994	0.934
HPC	0.830	0.632	0.654	0.5295	0.887	0.6355	0.906	0.643	0.947	0.927	0.194	0.936	0.924	0.909
Apache	1.000	0.000	1.000	0.694	1.000	0.694	0.314	0.0065	1.000	0.994	0.012	0.928	1.000	1.000
HealthApp	0.174	0.129	0.639	0.152	0.780	0.1085	0.279	0.112	1.000	0.6685	0.548	0.944	1.000	1.000
Mac	0.698	0.125	0.757	0.0325	0.787	0.218	0.520	0.169	0.778	0.490	0.548	0.771	0.831	0.82
Proxifier	0.508	0.000	0.527	0.000	0.527	0.000	0.027	0.000	1.000	0.000	0.025	0.895	1.000	0.999
Zookeeper	0.841	0.452	0.964	0.452	0.967	0.497	0.725	0.474	0.995	0.988	0.154	0.976	0.993	0.864
Thunderbird	0.943	0.026	0.844	0.027	0.955	0.047	0.189	0.004	0.257	0.473	0.256	0.971	0.991	0.933
Spark	0.884	0.004	0.905	0.3205	0.920	0.362	0.382	0.2585	0.4915	0.954	0.634	0.967	0.997	0.943
Android	0.880	0.714	0.919	0.245	0.911	0.709	0.791	0.413	0.885	0.331	0.523	0.842	0.983	0.965
Linux	0.701	0.122	0.605	0.088	0.690	0.184	0.147	0.124	0.389	0.388	0.185	0.971	0.937	0.868
Hadoop	0.885	0.0825	0.778	0.1125	0.948	0.269	0.428	0.113	0.787	0.384	0.291	0.949	0.989	0.902
OpenStack	0.743	0.019	0.764	0.000	0.733	0.019	0.236	0.000	0.503	0.872	0.092	0.744	0.957	0.938
Windows	0.566	0.1535	0.989	0.0035	0.997	0.159	0.695	0.1405	0.991	0.354	0.401	0.974	0.996	0.994
OpenSSH	0.925	0.133	0.554	0.1905	0.788	0.508	0.430	0.298	0.2295	0.9335	0.495	0.939	1.000	0.997
Average	0.766	0.167	0.792	0.208	0.865	0.319	0.478	0.18	0.723	0.62	0.309	0.920	0.975	0.942

617 3) *Baselines*: We compared the most advanced open-source
618 log parsing methods. LenMa [35] clusters logs based on log
619 similarity. Logram [9] distinguishes constant variables based
620 on the frequency of log tokens. Drain [12] clusters logs based
621 on rule trees. Spell [13] clusters logs based on the longest
622 identical subsequence between logs. LogPPT [15] uses 32 logs
623 to fine-tune the language model for log analysis. DivLog [17]
624 uses LLM for log parsing by adding contextual knowledge to
625 prompts.

626 B. Effectiveness of SelfLog

627 Table I displays the GA and PA of seven log parsing
628 methods across the 16 datasets. *SelfLog* outperformed the
629 other methods, achieving the highest average performance (see
630 the bottom line of Table I) in both GA and PA. *SelfLog* also
631 ranks as the best among existing algorithms in terms of
632 PTA and RTA. Due to space limitations, to compare with
633 more log parsers, we only selected the GA and PA to be
634 displayed in the table. The PTA and RTA of *SelfLog* are
635 shown in Table II below. It showed a 12.7% improvement
636 in GA compared to Drain and a 51.9% improvement in PA
637 compared to LogPPT. LogPPT and Logram methods are the
638 most unstable. The accuracy of Logram on Proxifier dataset
639 is only 0.027. This is because variables appear repeatedly
640 in Proxifier dataset, causing many variables to be incorrectly
641 recognized as constants. Drain has also achieved good results
642 in both stability and average GA, but in Proxifier dataset
643 the GA of Drain is only 0.527. Because all logs start with
644 variables, Drain needs to perform group analysis based on the
645 first few tokens, so the effect will be poor. This is because
646 Drain assumes that the initial tokens in logs are constants,
647 but in the Proxifier dataset, the majority of logs start with
648 variables, leading to misjudgments by Drain.

649 C. Efficient and Cost of SelfLog

650 As shown in Table I, DivLog is the best method apart
651 from *SelfLog*. Both our *SelfLog* and DivLog are based on
652 LLMs, and the two most important metrics for using LLMs are
653 processing time and cost. Therefore, we used 2000 logs from
654 five representative datasets to compare the processing time of
655 the two methods, as well as the number of input and output
656 tokens, since LLMs are billed based on the number of tokens.
657 In addition to these, we also detail PTA, RTA, PA, and GA as
658 accuracy criteria. From Fig. 6, it can be seen that *SelfLog* is
659 significantly lower than DivLog in both processing time and
660 token size. Under the circumstances that the log processing
661 accuracy of *SelfLog* is better than that of DivLog (as shown
662 in Fig. 6 (d), with a 190.5% improvement in the GA metric
663 and a 9.1% improvement in the PA metric), the processing
664 time of Proxifier dataset for *SelfLog* is only 1% of that for
665 DivLog, and the number of tokens is $\frac{1}{10}$ that of DivLog. The
666 main reason behind this is that DivLog requires a call to the
667 LLM for each log entry, whereas *SelfLog*, through N-Gram-
668 based Grouper and Log Hitter, only needs to call the LLM
669 when a new group appears. Since the LLM is currently the
670 bottleneck in log processing, reducing the number of calls to
671 the LLM can greatly improve the efficiency of log processing.
672 Moreover, as shown in Fig. 7(b), giving a group of logs to
673 the LLM for template extraction can better assist the model in
674 finding differences in the logs, thereby preparing to identify
675 constants and variables, and thus achieving better results.

676 D. Ablation Study

677 As shown in Table II, we sequentially removed the grouper,
678 parser, and merger of *SelfLog* to observe changes in the
679 model across four evaluation metrics. We did not perform
680 an ablation on Log Hitter because it does not contribute
681 to accuracy. Its main function is akin to a cache, capable
682 of storing previously parsed log groups and their templates

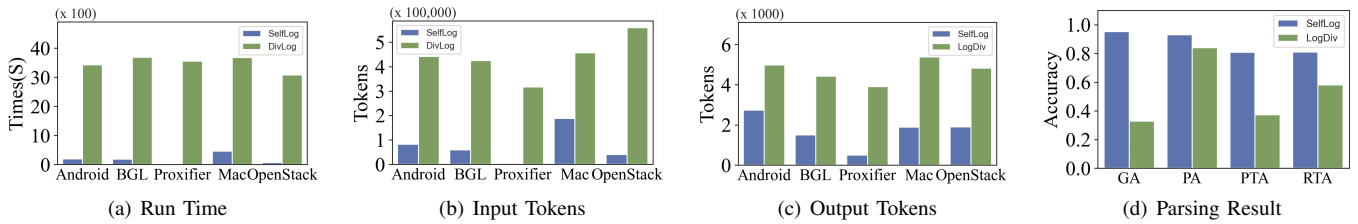


Fig. 6. Comparative histogram of log parsing effect between *SelfLog* and *DivLog* at running time, number of input tokens and output tokens when calling API, and log parsing effect.

683 for quick retrieval, eliminating the need for additional LLM
 684 invocations. The second row of Table II indicates that the
 685 component most affected within the entire *SelfLog* is the N-
 686 Gram-based Grouper. Upon its removal, GA dropped by 0.632,
 687 PTA by 0.658, RTA by 0.285, and PA by 0.19. Concurrently,
 688 the number of invocations of the LLM by *SelfLog* increased,
 689 leading to a significant rise in overall input and output tokens.
 690 The decline in efficiency is mainly due to the absence of
 691 grouping, every log entry requires an invocation of the LLM.
 692 The reason for the decline in effectiveness is illustrated in
 693 Fig. 7(b) with a detailed example, showing that presenting
 694 logs to the LLM in a grouped manner, as opposed to one by
 695 one, is more beneficial for template extraction, as the model
 696 can more accurately determine variables by comparing logs
 697 within the group.

698 Besides the Grouper, the second most impactful module on
 699 effectiveness is the LLM-based Log Parser, with declines of
 700 at least 0.5 in PA, PTA, and RTA. This is because, compared
 701 to statistical rule-based log parsing methods, LLM can make
 702 better judgments on whether each token is a variable or a
 703 constant leveraging its powerful natural language processing
 704 abilities. Without the LLM module, even with the presence
 705 of the Grouper, the accuracy of PA could only reach 0.434.
 706 Although the effectiveness of log parsing is already relatively
 707 high after the N-Gram-based Grouper and LLM-based Log
 708 Parser, Table II also shows that the final Tree-based Merger
 709 can enhance PA, PTA, and RTA one step further (more than 0.1).
 710 This is because logs are generally produced in a streaming
 711 manner, and it is possible that within a certain input window,
 712 a particular variable’s token may appear frequently (as shown
 713 in Fig. 5) and be mistakenly identified as a constant. The
 714 Merger, through the construction of a token tree, can correct
 715 these misidentified variables, thereby improving the model’s
 716 performance.

717 E. Parameter Sensitivity

718 In this section, we explore the impact of hyperparameters of
 719 our model on the outcome. There are three hyperparameters for
 720 the entire *SelfLog* system: the threshold used when dividing
 721 groups with N-Gram, the number of *Input Logs* from the
 722 same group fed into the prompt during log parsing with LLM,
 723 and the number of *Self-evolution Examples* selected from the
 724 prompt database. Their respective results are displayed in
 725 Table III, and Fig. 7(a) and Fig. 7(b). Firstly, we evaluate
 726 the impact of varying the N-Gram threshold in the Grouper

TABLE II
 ABLATION STUDY RESULTS OF *SelfLog*. THE LAST THREE LINES
 RESPECTIVELY REPRESENT THE PARSING EFFECT AFTER REMOVING
 DIFFERENT COMPONENTS FROM *SelfLog* .

Variants	GA	PA	PTA	RTA
<i>SelfLog</i>	0.975	0.942	0.876	0.873
- N-Gram-based Grouper	0.343	0.752	0.218	0.588
- LLM-based Log Parser	0.943	0.434	0.345	0.346
- Tree-based Merger	0.932	0.837	0.626	0.791

TABLE III
 THE AVERAGE GA UNDER DIFFERENT THRESHOLDS OF PILAR AND
SelfLog ON 16 DATASETS, THE IMPROVED EFFECT IS THE IMPROVEMENT
 OF *SelfLog* RELATIVE TO *DivLog*. *lines* REPRESENTS THE TOTAL
 NUMBER OF LOG ENTRIES.

Threshold	GA of PILAR	GA of <i>SelfLog</i>	Improved effect
threshold=0.10	0.81	0.876	8.14%
threshold=0.15	0.82	0.876	6.82%
threshold=0.20	0.82	0.877	6.95%
threshold=0.25	0.82	0.874	6.58%
threshold=0.30	0.79	0.870	10.12%
threshold=0.35	0.80	0.891	11.37%
threshold=0.40	0.81	0.891	10.00%
threshold=0.45	0.80	0.889	11.13%
threshold=0.50	0.79	0.889	12.53%
threshold=1/ <i>lines</i> * 5	-	0.877	6.95%
fluctuation range	0 ~ 0.03	0 ~ 0.019	-

727 on GA. We also examine the effects on other metrics such
 728 as PA, PTA, and RTA with parameter variation, with similar
 729 conclusions. Table III shows that our method maintains a high
 730 level of performance across different threshold values, with
 731 an improvement of at least 6.82% over *DivLog* [17], ranging
 732 from 0.876 to 0.891. Compared to PILAR [10], a method
 733 specifically optimized for parsing robustness, our fluctuation
 734 across different parameters is 0.019, which is 63% of PI-
 735 LAR’s fluctuation (0.019 v.s. 0.3), where a smaller fluctuation
 736 indicates better stability. It is noteworthy that the grouping
 737 threshold can be removed one step further. We propose a
 738 heuristic rule that the threshold for determining whether a
 739 token is variable using N-Gram can be dynamically adjusted
 740 by the total number of log lines, i.e., $\frac{1}{lines*5}$.

741 Regarding the number of representation logs in the same
 742 group for extracting the template, Fig. 7(a) reflects that the
 743 model stabilizes when the number of log entries exceeds 3.

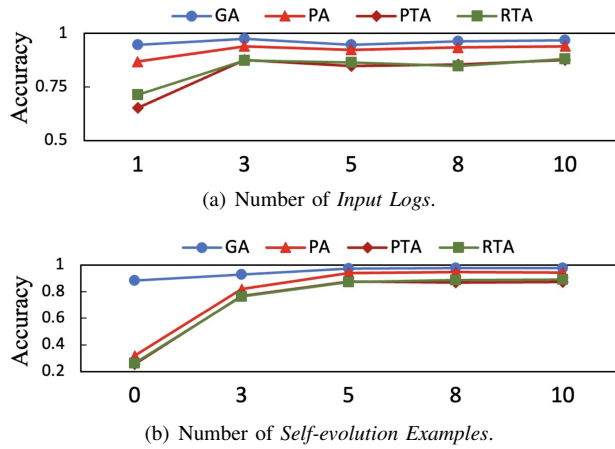


Fig. 7. The impact of varying quantities of *Input Logs* and *Self-Evolution Examples* on model performance.

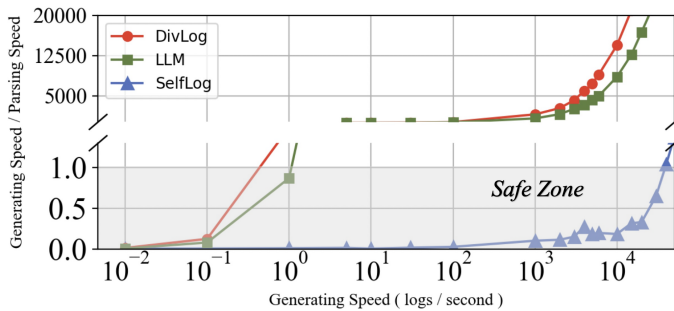


Fig. 8. Parsing speed of different LLM-based log parsing methods.

744 When the number of log entries for the same group increases
 745 from 1 to 3, GA improves by 6.7%, and PA by 8.2%, with
 746 the specific reasons introduced in Section III and Fig. 2 of
 747 the paper. As shown in Fig. 7(b), it is evident that without
 748 self-evolution examples, the model performs poorly. When
 749 the number of self-evolution examples increases from 0 to 3,
 750 PA improves significantly from 0.3 to 0.82. However, when
 751 the number of selected examples exceeds 5, the model’s
 752 performance tends to converge. This is because we use an
 753 Approximate Nearest Neighbors (ANN) method to select self-
 754 evolution examples from the prompt database, ensuring that
 755 as long as there are relevant logs, they can be retrieved.
 756 Thanks to LLM’s powerful few-shot learning capabilities, we
 757 can achieve good results with few relevant examples. Further
 758 adding examples yields marginal improvements.

759 F. Parsing Speed

760 While employing LLMs for log parsing offers numerous
 761 advantages, such as their strong semantic understanding
 762 capabilities and the ability of ICL to enhance the results of
 763 log parsing, the reality is that existing logs are typically
 764 generated in a streaming fashion and require real-time template
 765 extraction for immediate downstream anomaly detection. It
 766 is quite common for a large distributed system to generate
 767 tens of thousands of logs per second. However, existing
 768 algorithms such as DivLog are constrained by the generation

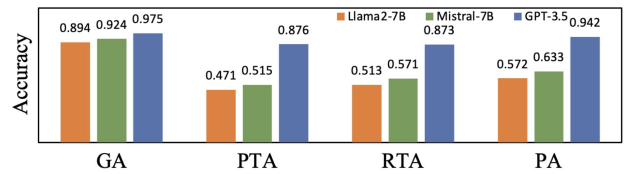


Fig. 9. A comparison of the performance of *SelfLog* when using different models as its backbone.

769 speed of LLM themselves. The generation speed of current
 770 large models is about 100 tokens per second [36], and a
 771 single log typically contains between 10 to 100 tokens in
 772 LogPAI [1], which means the rate can only reach a few logs
 773 per second. In contrast, our *SelfLog* benefits from a group-
 774 wise parsing paradigm and the caching mechanism of the log
 775 hitter, which significantly reduces the number of calls to the
 776 LLMs. As a result, the LLM is no longer a bottleneck. To
 777 get the exact parsing speed for existing LLM-based methods,
 778 including *SelfLog*, We use logs from HDFS [1] as input
 779 data, with 11,175,629 logs available. We replay these logs
 780 at different rates to test the log parsing speed of various
 781 models. In the experiment, we conduct multiple trials, each
 782 with a varying log generation speed, as shown in Fig. 8,
 783 where we test log generation speeds from 0.01 logs per second
 784 to 50,000 logs per second. We monitor the processing speeds
 785 of DivLog [17], vanilla LLM, and *SelfLog*, calculating the ratio
 786 of log generation speed to log parsing speed as the Y-axis.
 787 A ratio of less than 1.0 indicates that the log parsing speed
 788 exceeds the log generation speed, suggesting the model has
 789 sufficient capacity to handle more logs. Conversely, a ratio
 790 greater than 1 means the log parsing speed is less than the
 791 log generation speed, leading to a continuous backlog and,
 792 over time, potential Out-Of-Memory (OOM) issues. In Fig. 8,
 793 we mark the area where the Y-axis is less than 1 as the “safe
 794 zone”. When Y equals 1, the corresponding X value represents
 795 the peak parsing speed supported by the algorithm. It can be
 796 observed from the figure that existing LLM-based log parsing
 797 algorithms, which require an LLM call for each log, have
 798 processing speeds of fewer than 10 logs per second and are
 799 already beyond the “safe zone” when the log generating speed
 800 exceeds 10 logs per second, resulting in a backlog. In contrast,
 801 *SelfLog* remains within the “safe zone” even when the log rate
 802 is 10,000 per second and reaches a remarkable peak parsing
 803 speed of 45,000 logs per second.

804 G. Model backbone

805 Fig. 9 demonstrates the performance of *SelfLog* when uti-
 806 lizing different LLMs as the backbone. It is evident that
 807 as the capabilities of the LLMs improve, the performance
 808 of *SelfLog* also continuously enhances. Due to our resource
 809 limitations, we have only tested the 7-billion-parameter open-
 810 source model. We believe that with the ongoing advancement
 811 of the LLM community, *SelfLog* can achieve further improve-
 812 ments in the future.

V. THREATS TO VALIDITY

External Validity: In this article, we study and compare the effects of *SelfLog* and six state-of-the-art log parsing algorithms on 16 open-source datasets of LogPai [1]. Although these 16 datasets come from different systems, each dataset only has 2,000 manually labeled data, which does not represent logs in real scenarios. In the future, more realistic hand-labeled log datasets can be constructed to optimize the evaluation of various log parsers. We tested the efficiency and effect of *SelfLog* when processing a large number of logs in online work. Only testing the HDFS dataset cannot comprehensively and accurately display the online work efficiency of *SelfLog*. Further testing in real scenarios is needed.

Internal Validity: In the future, with the improvement of model capabilities, N-Gram-based Grouper may become a bottleneck limiting the effect of LLM on log analysis. When there is an error in classifying logs belonging to different templates into the same group, it will directly affect the final parsing results. But currently, *SelfLog* is still a robust, effective, and efficient log parsing algorithm.

Construct Validity: We set the *temperature* parameter of LLM as 0 to reduce the randomness of the results returned by LLM, but the results returned by LLM for the same input are still inconsistent. We record the experimental results through multiple experiments. Though ANN is better than the KNN used by DivLog [17] in terms of efficiency, it is not as good as KNN (K-Nearest Neighbors) in terms of retrieval accuracy.

VI. RELATED WORKS

A. Unsupervised log parsers

Unsupervised log parser does not require manual annotation of data for training and can be directly used in different systems for log parsing. Unsupervised log parsers can be further divided into frequent pattern mining-based [9]–[11], clustering-based [37]–[39], heuristic rule-based [12]–[14], and LLM-based methods [17]. Methods based on frequent pattern mining start from the data distribution itself and rely on data features (e.g. token frequent) to propose templates. The advantage is that it doesn't rely on artificially designed hyperparameters based on the data itself, and the method is highly robust (PILAR [10]). The disadvantage is that it is easily affected by the imbalance of data distribution. Logram [9] and LogCluster [11] all perform log analysis by extracting frequent patterns from logs. The clustering-based method adopts grouping first. By default, logs in the same group have the same template. Templates are proposed based on the differences in logs in the same group (different tokens are replaced with $\langle * \rangle$). LogMine [37] and LogTree [38] use the hierarchical clustering method to group logs, and LTE [39] use density-based clustering to group logs. LenMa [35] and FLP [40] adopt online grouping strategies to support online parsing. Based on the heuristic rule method, human knowledge is transformed into rules for log analysis by carefully observing the data. Drain [12], Spell [13], and IPLoM [14] have achieved good log parsing results by fine-tuning algorithm hyperparameters for different data.

However, due to algorithm design flaws, they cannot correctly parse all log types and have poor robustness. The LLM-based method directly utilizes LLM's powerful natural language understanding capabilities. By providing a few context examples to build prompts, DivLog [17] has achieved the most advanced results in PA.

B. Supervised log parsers

Supervised log parsers usually use deep learning methods to train or fine-tune models by manually annotating data. VALB [25] manually annotate constants and variable categories using a method similar to named entity recognition, using the BiLSTM [26] model to understand and perform template extraction and variable category annotation. SemParser [16] extracts concept-instance (CI) pairs through the designed semantic miner, and then uses the joint parser to combine the context information to identify variables. LogPPT [15] proposes to use a small number of logs and template examples to fine-tune the pre-trained model RoBERTa [41] and then perform log analysis. However, the computation cost of fine-tuning is negligible.

VII. CONCLUSION

The advent of LLMs has presented a promising alternative for accurate log parsing, yet they come with their own set of challenges, particularly the need for manual annotation and the inefficiency of processing large volumes of logs. To overcome these obstacles, we introduce *SelfLog*, a self-evolving log parsing method that leverages the power of LLMs while mitigating their limitations. Our approach operates in two innovative ways: firstly, by using similar history $\langle \text{group}, \text{template} \rangle$ pairs outputted by LLM itself, which serves as prompts for new log entries, thus allowing the model to evolve and learn autonomously without the need for manual labeling. Secondly, we implement an N-Gram-based grouper and log hitter mechanism, which enhances the parsing performance by processing logs in groups rather than individually and significantly reduces redundant calls to LLMs for logs whose group templates have been previously extracted. Our comprehensive evaluation across 16 public datasets, encompassing tens of millions of logs, has demonstrated that *SelfLog* not only achieves state-of-the-art performance with a GA of 0.984 and a PA of 0.743 but also excels in efficiency, processing at a remarkable speed (over 45,000 logs per second) compared with existing LLM-based log parsing methods. In a nutshell, by integrating N-Gram-based grouping with self-evolutionary in-context learning, *SelfLog* fully harnesses the advantages of LLM in few-shot learning while avoiding inefficiency pitfalls. We will continue to explore the application of this paradigm in log analysis in the future.

REFERENCES

- [1] J. Zhu, S. He, J. Liu, P. He, Q. Xie, Z. Zheng, and M. R. Lyu, "Tools and benchmarks for automated log parsing," in *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 2019, pp. 121–130.

- 920 [2] Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and
921 G. Batista, "The ucr time series classification archive," July 2015, www.
922 cs.ucr.edu/~eamonn/time_series_data/.
- 923 [3] D. Zhang, X. Peng, C. Sha, K. Zhang, Z. Fu, X. Wu, Q. Lin, and
924 C. Zhang, "Deeptralog: Trace-log combined microservice anomaly de-
925 tection through graph-based deep learning," in *Proceedings of the 44th*
926 *International Conference on Software Engineering*, 2022, pp. 623–634.
- 927 [4] S. He, J. Zhu, P. He, and M. R. Lyu, "Experience report: System
928 log analysis for anomaly detection," in *2016 IEEE 27th international*
929 *symposium on software reliability engineering (ISSRE)*. IEEE, 2016,
930 pp. 207–218.
- 931 [5] M. Du, F. Li, G. Zheng, and V. Srikumar, "Deeplog: Anomaly detection
932 and diagnosis from system logs through deep learning," in *Proceedings*
933 *of the 2017 ACM SIGSAC conference on computer and communications*
934 *security*, 2017, pp. 1285–1298.
- 935 [6] J. Liu, J. Zhu, S. He, P. He, Z. Zheng, and M. R. Lyu, "Logzip: Extract-
936 ing hidden structures via iterative clustering for log compression," in
937 *2019 34th IEEE/ACM International Conference on Automated Software*
938 *Engineering (ASE)*. IEEE, 2019, pp. 863–873.
- 939 [7] R. Tian, Z. Diao, H. Jiang, and G. Xie, "Logdac: A universal efficient
940 parser-based log compression approach," in *ICC 2022-IEEE Interna-*
941 *tional Conference on Communications*. IEEE, 2022, pp. 3679–3684.
- 942 [8] S. Locke, H. Li, T.-H. P. Chen, W. Shang, and W. Liu, "Logassist:
943 Assisting log analysis through log summarization," *IEEE Transactions*
944 *on Software Engineering*, vol. 48, no. 9, pp. 3227–3241, 2021.
- 945 [9] H. Dai, H. Li, C.-S. Chen, W. Shang, and T.-H. Chen, "Logram: Efficient
946 log parsing using n n-gram dictionaries," *IEEE Transactions on Software*
947 *Engineering*, vol. 48, no. 3, pp. 879–892, 2020.
- 948 [10] H. Dai, Y. Tang, H. Li, and W. Shang, "Pilar: Studying and mitigating
949 the influence of configurations on log parsing," in *2023 IEEE/ACM 45th*
950 *International Conference on Software Engineering (ICSE)*. IEEE, 2023,
951 pp. 818–829.
- 952 [11] R. Vaarandi and M. Pihelgas, "Logcluster-a data clustering and pattern
953 mining algorithm for event logs," in *2015 11th International conference*
954 *on network and service management (CNSM)*. IEEE, 2015, pp. 1–7.
- 955 [12] P. He, J. Zhu, Z. Zheng, and M. R. Lyu, "Drain: An online log parsing
956 approach with fixed depth tree," in *2017 IEEE international conference*
957 *on web services (ICWS)*. IEEE, 2017, pp. 33–40.
- 958 [13] M. Du and F. Li, "Spell: Streaming parsing of system event logs," in
959 *2016 IEEE 16th International Conference on Data Mining (ICDM)*.
960 IEEE, 2016, pp. 859–864.
- 961 [14] A. Makanju, A. N. Zincir-Heywood, and E. E. Milios, "A lightweight
962 algorithm for message type extraction in system application logs," *IEEE*
963 *Transactions on Knowledge and Data Engineering*, vol. 24, no. 11, pp.
964 1921–1936, 2011.
- 965 [15] V.-H. Le and H. Zhang, "Log parsing with prompt-based few-shot
966 learning," in *2023 IEEE/ACM 45th International Conference on Software*
967 *Engineering (ICSE)*. IEEE, 2023, pp. 2438–2449.
- 968 [16] Y. Huo, Y. Su, C. Lee, and M. R. Lyu, "Semparser: A semantic parser
969 for log analytics," in *2023 IEEE/ACM 45th International Conference on*
970 *Software Engineering (ICSE)*. IEEE, 2023, pp. 881–893.
- 971 [17] J. Xu, R. Yang, Y. Huo, C. Zhang, and P. He, "Divlog: Log parsing with
972 prompt enhanced in-context learning," in *Proceedings of the IEEE/ACM*
973 *46th International Conference on Software Engineering*, 2024, pp. 1–12.
- 974 [18] X. Zhang, Y. Xu, Q. Lin, B. Qiao, H. Zhang, Y. Dang, C. Xie,
975 X. Yang, Q. Cheng, Z. Li *et al.*, "Robust log-based anomaly detection on
976 unstable log data," in *Proceedings of the 2019 27th ACM joint meeting*
977 *on European software engineering conference and symposium on the*
978 *foundations of software engineering*, 2019, pp. 807–817.
- 979 [19] X. Wang, X. Zhang, L. Li, S. He, H. Zhang, Y. Liu, L. Zheng, Y. Kang,
980 Q. Lin, Y. Dang *et al.*, "Spine: A scalable log parser with feedback
981 guidance," in *Proceedings of the 30th ACM Joint European Software*
982 *Engineering Conference and Symposium on the Foundations of Software*
983 *Engineering*, 2022, pp. 1198–1208.
- 984 [20] S. He, P. He, Z. Chen, T. Yang, Y. Su, and M. R. Lyu, "A survey
985 on automated log analysis for reliability engineering," *ACM computing*
986 *surveys (CSUR)*, vol. 54, no. 6, pp. 1–37, 2021.
- 987 [21] Y. Gan, M. Liang, S. Dev, D. Lo, and C. Delimitrou, "Sage: practical
988 and scalable ml-driven performance debugging in microservices," in
989 *Proceedings of the 26th ACM International Conference on Architectural*
990 *Support for Programming Languages and Operating Systems*, 2021, pp.
991 135–151.
- 992 [22] H. Wang, Z. Wu, H. Jiang, Y. Huang, J. Wang, S. Kopru, and T. Xie,
993 "Groot: An event-graph-based approach for root cause analysis in
994 industrial settings," in *2021 36th IEEE/ACM International Conference*
995 *on Automated Software Engineering (ASE)*. IEEE, 2021, pp. 419–429.
- [23] D. Yuan, H. Mai, W. Xiong, L. Tan, Y. Zhou, and S. Pasupathy, "Sherlog:
996 error diagnosis by connecting clues from run-time logs," in *Proceedings*
997 *of the fifteenth International Conference on Architectural support for*
998 *programming languages and operating systems*, 2010, pp. 143–154.
- [24] W. Xu, L. Huang, A. Fox, D. Patterson, and M. I. Jordan, "Detecting
1000 large-scale system problems by mining console logs," in *Proceedings*
1001 *of the ACM SIGOPS 22nd symposium on Operating systems principles*,
1002 2009, pp. 117–132.
- [25] Z. Li, C. Luo, T.-H. Chen, W. Shang, S. He, Q. Lin, and D. Zhang, "Did
1004 we miss something important? studying and exploring variable-aware
1005 log abstraction," in *2023 IEEE/ACM 45th International Conference on*
1006 *Software Engineering (ICSE)*. IEEE, 2023, pp. 830–842.
- [26] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with
1008 deep recurrent neural networks," in *2013 IEEE international conference*
1009 *on acoustics, speech and signal processing*. Ieee, 2013, pp. 6645–6649.
- [27] S. C. Fanni, M. Febi, G. Aghakhanyan, and E. Neri, "Natural language
1011 processing," in *Introduction to Artificial Intelligence*. Springer, 2023,
1012 pp. 87–99.
- [28] N. Wies, Y. Levine, and A. Shashua, "The learnability of in-context
1014 learning," *Advances in Neural Information Processing Systems*, vol. 36,
1015 2024.
- [29] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*,
1017 "Language models are unsupervised multitask learners," *OpenAI blog*,
1018 vol. 1, no. 8, p. 9, 2019.
- [30] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun,
1020 J. Xu, and Z. Sui, "A survey on in-context learning," *arXiv preprint*
1021 *arXiv:2301.00234*, 2022.
- [31] S. J. Reddi, S. Miryoosefi, S. Karp, S. Krishnan, S. Kale, S. Kim,
1023 and S. Kumar, "Efficient training of language models using few-shot
1024 learning," in *International Conference on Machine Learning*. PMLR,
1025 2023, pp. 14553–14568.
- [32] Z. Wang, W. Zhong, Y. Wang, Q. Zhu, F. Mi, B. Wang, L. Shang,
1027 X. Jiang, and Q. Liu, "Data management for large language models: A
1028 survey," *arXiv preprint arXiv:2312.01700*, 2023.
- [33] G. A. Miller, "Wordnet: a lexical database for english," *Communications*
1030 *of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [34] Z. A. Khan, D. Shin, D. Bianculli, and L. Briand, "Guidelines for
1032 assessing the accuracy of log message template identification techni-
1033 ques," in *Proceedings of the 44th International Conference on Software*
1034 *Engineering*, 2022, pp. 1095–1106.
- [35] K. Shima, "Length matters: Clustering system log messages using length
1036 of words," *arXiv preprint arXiv:1611.03213*, 2016.
- [36] S. Kou, L. Hu, Z. He, Z. Deng, and H. Zhang, "Cilms: Consistency
1038 large language models," *arXiv preprint arXiv:2403.00835*, 2024.
- [37] H. Hamooni, B. Debnath, J. Xu, H. Zhang, G. Jiang, and A. Mueen,
1040 "Logmine: Fast pattern recognition for log analytics," in *Proceedings*
1041 *of the 25th ACM international on conference on information and*
1042 *knowledge management*, 2016, pp. 1573–1582.
- [38] L. Tang and T. Li, "Logtree: A framework for generating system events
1044 from raw textual logs," in *2010 IEEE International Conference on Data*
1045 *Mining*. IEEE, 2010, pp. 491–500.
- [39] J. Ya, T. Liu, H. Zhang, J. Shi, and L. Guo, "An automatic approach to
1047 extract the formats of network and security log messages," in *MILCOM*
1048 *2015-2015 IEEE Military Communications Conference*. IEEE, 2015,
1049 pp. 1542–1547.
- [40] Y. Zhong, Y. Guo, and C. Liu, "Flp: a feature-based method for log
1051 parsing," *Electronics letters*, vol. 54, no. 23, pp. 1334–1336, 2018.
- [41] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis,
1053 L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert
1054 pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.
- 1055