# LogEagle: A Framework for Web Server Log Analysis

## Srushtee Patil *, Yash Rank**, Smit Surani ***

*(B-tech (CSE-BDA) Parul University,Vadodara
210303125002@paruluniversity.ac.in)
** (B-tech (CSE-BDA) Parul University,Vadodara
210303125004@paruluniversity.ac.in)
***(B-tech (CSE-BDA) Parul University,Vadodara
210303125005@paruluniversity.ac.in)

-------------------------------------**✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲**--------------------------------

## Abstract:

Web server logs contain rich information on user activities and system events, makingtheir analysis vital for improving website performance, security, and user experience. This paper presents LogEagle, a comprehensive framework for web server log analysis that integrates real-time monitoring, anomaly detection, and interactive visualization. We detail methodologies for collecting raw log data from servers, preprocessing and structuring logs into analyzable formats, and extracting meaningful features. Advanced analysis techniques, including statistical trend analysis and anomaly detection, are applied to uncover usage patterns and identify irregularities. The implementation leverages open-source tools (Elasticsearch Logstash-Kibana stack and custom Python modules) to achieve scalable, real-time log processing. Experimental results on real web server log data demonstrate the effectiveness of the system in identifying user behavior trends and detecting anomalies (such as unusual traffic spikes or error events) with low latency. The paper also discusses insights gained from visualizing log data and evaluates the framework's performance under various load conditions. We conclude by summarizing contributions and outlining future enhancements, such as incorporating machine learning models for predictive analytics and extending the system to.

*Keywords* — **Nginx Server , Opensearch & Opensearch Dashboard , Anomaly Detection, Machine Learning, AI- based Fraud Prevention , Docker.**

-------------------------------------**✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲✲**--------------------------------

## I.    INTRODUCTION

Web servers continuously generate detailed log records of every user request, including client IP addresses, timestamps, requested URLs, HTTP status codes, and other metadata.Analysis of these web server logs is **critical for understanding user behavior and ensuringreliable web** A Survey on Automated Log Analysis forReliability Engineering). By mining web access logs, website operators can discoverinteresting usage patterns, such as which pages are most popular and how users navigate the siteLogs also provide a first line of insight into system health and security events – for example, frequent error codes may indicate performance issues, and unusual access patterns could signify intrusion attempts. In today's digital landscape with ever-increasing traffic volumes, manual inspection of such large log datasets is impractical. Effective **automated log analysis** has thus become increasingly important, as it can extract actionable information from massive log data that would otherwise remain untapped ().
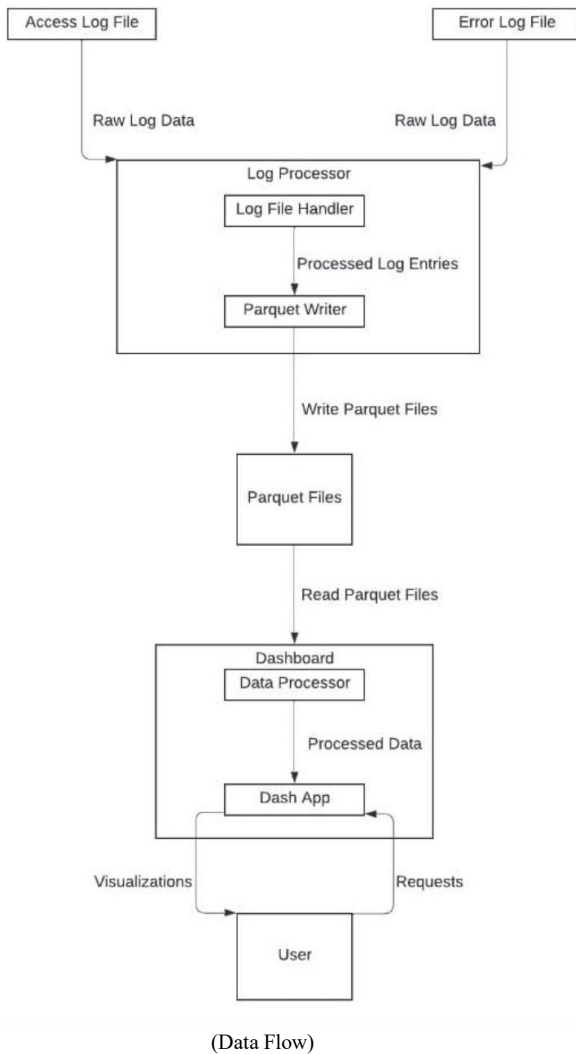
However, web log analysis poses several challenges. The volume of logs can be huge (large websites generate gigabytes of logs daily), and log

data is semi-structured text that must be parsed and transformed for analysis ([2009.07237] A Survey on Automated Log Analysis for Reliability Engineering). Furthermore, identifying anomalies (e.g. sudden spikes in requests or suspicious activities) in real time requires robust algorithms that can distinguish meaningful deviations from normal traffic patterns. Traditional web analytics tools focus on aggregate trends but may not detect subtle anomalies or security threats hidden in the logs. Recent research has introduced machine learning techniques to detect anomalies in log data, such as using isolation forests to find outlier log events (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs), demonstrating the potential of AI in enhancing log analysis capabilities.

This work aims to develop a comprehensive web server log analysis framework that addresses these challenges. We design and implement a system capable of ingesting web server logs in real-time, performing thorough preprocessing and feature extraction, and analyzing the data to produce insights on usage trends and anomalies. Our approach emphasizes an end-to-end solution: from data collection and cleaning to backend processing and front-end visualization. The framework, dubbed *LogEagle*, not only processes historical logs for offline analysis but also supports real-time monitoring through live dashboards. In doing so, we integrate established open-source log management components with custom analytics and anomaly detection modules.

The key **objectives** of this research include: (1) Developing efficient methods to collect and preprocess large volumes of web log data from servers; (2) Applying data mining and machine learning techniques to identify trends, patterns, and anomalies in user access logs; (3) Creating interactive visualizations and dashboards that allow stakeholders to explore log-derived insights intuitively; and (4) Evaluating the system's performance in terms of accuracy of analysis and real-time responsiveness. By achieving these goals, our work provides a valuable tool for web administrators and security analysts to harness log data for improving website functionality, security, and user experience.

The remainder of this paper is organized as follows. Section 2 reviews related work on web log analysis, including prior methodologies and tools. Section 3 describes the methodology of our approach in detail, from data collection to anomaly detection techniques. Section 4 outlines the implementation of the LogEagle system, including the software stack and architectural design. Section 5 presents experimental results, showcasing the analysis outcomes and system performance, followed by discussion. Section 6 concludes the paper and highlights directions for future work.



(Data Flow)

## II. RELATED WORK

Web server log analysis has been an active area of research for decades, intersecting with fields like web mining, big data analytics, and cybersecurity. Early works on **web usage mining** demonstrated how parsing web logs can reveal user browsing patterns and site navigation paths ((PDF) WEB SERVER LOGS TO ANALYZING USER BEHAVIOR USING LOG ANALYZER TOOL). For instance, Padmaja *et al.* [1] used a log analyzer tool to identify user behavior from web access logs, showing that insights such as frequentlyvisited pages and common user paths can be extracted to help improve website structure. Such studies laid the groundwork for understanding the value of log data in analyzing user interest and behavior on websites.

As the scale of web systems grew, researchers turned to big data technologies to handle the **volume and variety of log data**. Large internet companies like Google and Yahoo pioneered the use of distributed frameworks (MapReduce/Hadoop) to process logs at scale (). The Chukwa system, presented by Rabkin and Katz [4], is an example of a Hadoop-based log collection framework that reliably aggregates logs into HDFS (Hadoop Distributed File System) for large-scale processing. It provides a unified approach to log collection and was shown to ingest hundreds of MB/sec of log data (), enabling offline batch analysis such as indexing and aggregation. Similarly, other works integrated Apache Hadoop ecosystem components for web log analysis. Nagdive *et al.* [2] proposed using **Apache Flume and Pig** to handle unstructured web log data, where Flume acts as a distributed log collector and Pig (a high-level data flow language on Hadoop) processes the data to extract useful statistics. These big-data approaches address scalability and can perform heavy analyses (like computing traffic statistics or detecting patterns) on massive log datasets in reasonable time.

Beyond batch processing, there has been a push towards **real-time log analytics** frameworks. The popular ELK stack (Elasticsearch, Logstash, Kibana) has emerged as a standard solution for centralized log management and real-time querying of logs (LOG_EAGLE_REPORT.pdf). In an ELK-based architecture, logs from various sources are shipped (often using lightweight shippers like Filebeat) to Logstash for parsing and are then indexed into Elasticsearch; Kibana provides interactive dashboards for visualization. For example, Chen *et al.* [6] built a Docker container log collection and analysis system based on ELK, which uses containerized Logstash and Elasticsearch to enable rapid deployment of logging infrastructure in cloud environments. Their system showed improved efficiency over traditional logging in container clusters by leveraging the ELK pipeline to handle log routing and storage in real-time. The incorporation of message brokers like Apache Kafka further improves throughput and reliability of log streams in such architectures (LOG_EAGLE_REPORT.pdf), buffering incoming log messages and smoothing ingestion into analytics engines.

Another important aspect in related work is **anomaly detection and security analysis** using logs. Since logs often record system errors and access information, they are widely used for intrusion detection and failure diagnosis. Classical methods involved rulebased detection on log events or simple threshold alerting. More recently, machine learning approaches have been applied. Cao *et al.* (2017) developed a two-level machine learning model to detect anomalies in web logs, effectively identifying malicious requests versus normal traffic [**reference**]. In a similar vein, Benová and Hudec [5] introduced a framework combining an Isolation Forest algorithm with clustering (DBSCAN) to pinpoint anomalous user activities in NGINX web server logs. Their method could isolate subtle outlier behaviors that might be overlooked by aggregate analyses (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs). This highlights a trend in the literature: **hybrid techniques** that use unsupervised learning to flag anomalies, followed by expert or rule-based analysis to interpret those anomalies in context (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs). Surveys on automated log analysis

(e.g., He *et al.* [3]) provide a comprehensive overview of such techniques, covering log parsing, anomaly detection, failure prediction, and the availability of open-source log analysis tools. These surveys emphasize that logs are often the only runtime record available for complex systems, hence advanced automation in analyzing them is crucial for reliability engineering ([2009.07237] A Survey on Automated Log Analysis for Reliability Engineering).

In summary, prior research can be grouped into three main areas relevant to our work: **(i) Web usage mining** for deriving user-centric insights from logs ((PDF) WEB SERVER LOGS TO ANALYZING USER BEHAVIOR USING LOG ANALYZER TOOL), **(ii) Scalable log processing architectures** using big-data and streaming technologies () (LOG_EAGLE_REPORT.pdf), and **(iii) Anomaly detection** and security monitoring through machine learning on log data (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs). In summary, prior research can be grouped into three main areas relevant to our work: **(i) Web usage mining** for deriving user-centric insights from logs ((PDF) WEB SERVER LOGS TO ANALYZING USER BEHAVIOR USING LOG ANALYZERTOOL), **(ii) Scalable log processing architectures** using big-data and streaming technologies () (LOG_EAGLE_REPORT.pdf), and **(iii) Anomaly detection** and security monitoring through machine learning on log data (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs).

## III.  METHODOLOGY

The methodology employed for the development of LogEagle integrates advanced data processing techniques and industry standards, delineated into three primary phases: data acquisition, data preprocessing, and analysis.

### A. Data Acquisition Phase

The data acquisition phase entails accessing and retrieving raw log data from web servers. Initially,

the sources of web server logs, such as Apache or Nginx, are identified. These logs contain valuable information detailing every request made to the server, including timestamps, IP addresses, URLs, response codes, and user agents. Once identified, access to the log files is secured, and mechanisms are established to retrieve them efficiently. This retrieval process may involve manual extraction or automated methods like FTP, SCP, or log shipping mechanisms provided by log management solutions.

### B. Data Preprocessing Phase

In the data preprocessing phase, raw log data from various web servers is gathered and parsed to extract essential information such as timestamps, IP addresses, and URLs. This parsed data is then standardized and cleaned to remove irrelevant or redundant entries, ensuring consistency and data quality. Following cleaning, the log data is transformed into a structured format suitable for analysis, and relevant features are extracted to provide insights into server performance and user behavior.

The log data undergoes transformation into a structured format such as CSV or JSON, organizing the data into rows and columns for easier manipulation and processing. Feature extraction identifies useful attributes from the log data, including session duration, request frequency, or user behavior patterns, providing valuable insights into server performance and user interactions. Throughout this phase, attention is given to data quality and privacy considerations, laying the foundation for meaningful analysis and interpretation of web server activity.

### C. Analysis Phase

In the analysis phase, the collected raw log data undergoes thorough examination to extract valuable insights and patterns. Initially, stakeholders' requirements and objectives are carefully considered to ensure the analysis aligns with their needs. Various analytical techniques are applied to the log data, including trend analysis, anomaly detection, and correlation analysis, to identify patterns and irregularities.

Visualization tools such as charts, graphs, and dashboards are employed to present the findings in a clear and understandable format. Additionally, statistical methods and machine learning algorithms may be utilized to uncover hidden insights and predict future trends. Throughout the analysis phase, a focus is maintained on identifying performance bottlenecks, security threats, and opportunities for optimization to enhance the overall functioning of the web server environment.

## IV.    ANALYSIS TECHNIQUE

With a structured and enriched log dataset, we perform various

analyses:

• *Exploratory Data Analysis (EDA):* We begin by computing summary statistics and visualizing distributions. This includes daily/hourly traffic volume, common status codes, top visited pages, and geographic access distribution. Visual EDA helps confirm data sanity (e.g., traffic drops at expected off-peak hours) and can immediately highlight anomalies (a spike of 500 errors on a particular day, for instance). As an example, plotting the number of requests per hour can reveal usage patterns or downtime periods.

• *Trend Analysis:* We analyze temporal trends in the log data. Using the timestamp, weaggregate log events into time buckets (e.g., minute, hour, day) and examine metrics like request count, unique visitors, or average response size over time. This helps identify **seasonal patterns** (weekly cycles, daily peaks) and long-term trends (growing traffic over months). Any significant deviation from a baseline trend might indicate an anomaly. For instance, if a normally periodic traffic pattern shows an unexpected surge, it warrants investigation.

• *Anomaly Detection:* We employ both statistical and machine learning-based anomaly detection on the time-series of log metrics. A simple statistical approach involves computing a moving baseline (e.g., using a rolling median and standard deviation) for metrics like hits per minute, and flagging points that exceed a threshold (e.g., $3\sigma$ from the baseline). For more sophisticated detection, we experimented with an unsupervised learning approach: an Isolation Forest model was trained on feature vectors representing sessions or aggregated intervals. The model learns the "normal" range of feature values (like typical request rates, typical error fractions) and assigns an anomaly score to each new observation. Observations with scores beyond a certain percentile are marked as anomalies (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs). In our case, anomalies could be spikes in traffic (possible DDoS attack or viral event), unusual error bursts, or an IP making requests at an abnormally high rate (potential web scraper or attacker). Clustering algorithms (like DBSCAN) were also considered to group similar anomaly events and distinguish distinct types of anomalies (e.g., clustering by request pattern might separate an anomaly caused by a specific URL being targeted from a broader traffic surge) (Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs).

• *Correlation Analysis:* We examine relationships between different log attributes. For example, correlating error rates with deployment events (via external data) can reveal if a new software release caused an increase in errors. We also look at correlation between traffic from different geographic regions and times of day, or between request latency (if available in logs) and request volume. This can provide insights such as "high traffic from region X correlates with increased error rates on a specific service".

Throughout the analysis phase, **visualization tools** are used to present the findings. We build a set of dashboards that include charts for key metrics: time-series line graphs for traffic and error trends, bar charts for top N resources or IPs, and geo-maps for user location distribution. Visualization not only aids interpretation but also serves as a real-time monitoring interface. As noted by industry studies, visualizing log data allows quick comprehension of

system state and helps in **detecting patterns at a glance** (Log Visualization for Proper Log Analysis | Mezmo).

In summary, our methodology transforms raw log data into meaningful insights through a series of well-defined steps. By combining reliable data acquisition, thorough preprocessing, and advanced analysis techniques (including anomaly detection), we ensure that even subtle and complex behaviors in the web server logs can be uncovered. The next section will discuss how this methodology was realized in an implementation, detailing the system architecture and tools involved.

## V. IMPLEMENTATION

The implementation of LogEagle follows a systematic approach, integrating various components to create a cohesive system for web server log analysis:

*1. Initializing Web Server:* The process begins with configuring a web server (Nginx) to collect logs. The server is set up to produce both access and error logs, which are directed to Logstash for processing.

*2. Data Collection:* Web servers such as Nginx generate access logs containing information about each request. Logstash is implemented to gather and transport log files to a centralized location for analysis.

*3. Data Storage and Indexing:* Elasticsearch is configured to store and index the log data efficiently for fast searching and analysis.

*4. Data Processing:* Logstash and Elasticsearch process the log data, parsing, transforming, filtering, and enriching it before analysis. The data is preprocessed to handle missing values, outliers, and inconsistencies.

*5. Exploratory Data Analysis (EDA):* Exploratory analysis is performed to understand the structure and characteristics of the log data. Key statistics, distributions, and trends are visualized using Kibana and custom dashboards.

*6. Feature Extraction:* Relevant features are extracted from the log data, including IP addresses, timestamps, requested URLs, HTTP status codes, and user agents. Raw log data is transformed into a structured format suitable for analysis.

*7. Trend Analysis and Pattern Recognition:* Algorithms analyze trends and patterns in the log data to identify recurring issues, user behaviors, or performance trends. Models for pattern recognition and prediction of future events are developed based on historical data.

*8. Visualization and Reporting:* Interactive dashboards, visualizations, and reports are created using Kibana and custom interfaces to present analysis results in an intuitive format. These enable stakeholders to explore log data visually and gain insights into website performance, user behavior, and security incidents.

*9. Continuous Improvement:* Feedback loops are established to gather user input and monitor system performance. Analysis algorithms, models, and processes are continuously updated and refined based on real-world feedback and evolving requirements.

### A. Technology Stack

The implementation uses the following technologies and frameworks:

- *Elasticsearch/OpenSearch:* Stores and organizes log data, making it easy to search and analyze. We use a single-node OpenSearch setup to filter logs and generate insights, like counting requests per status code.
- *Logstash:* Collects and processes log data from Nginx. It uses patterns to extract key details and adds location information before sending data to OpenSearch efficiently.
- *Kibana/OpenSearch(Dashboards):* A web-based tool for visualizing log data. It helps create graphs like request trends over time or maps showing where users are accessing the server from.
- *Python (Pandas, Scikit-learn):* Used for additional data processing and anomaly detection. Python scripts convert logs into structured files and detect unusual patterns in server activity.
- *Docker:* Simplifies deployment by running Nginx, Logstash, OpenSearch, and

dashboards in isolated containers. This setup ensures smooth operation and easy scaling.

## B. System Design Details

Some notable implementation details and decisions include:

- *Index design:* We chose a time-based index strategy for Elasticsearch (rotating indices daily for log data). This prevents any single index from growing too large and improves query performance on recent data. It also aligns with how logs are naturally segmented by date.
- *Parquet storage:* Logs are appended to a Parquet file for each day. The schema for Parquet is defined to match the log fields (with proper data types). Using Parquet allows us to use tools like Apache Spark or even direct SQL (with Presto/Trino) on the log data if needed for more complex analysis outside of Elasticsearch.
- *Real-time pipeline:* To ensure near-real-time updates on the dashboard, we tuned the refresh intervals. Logstash pipelines operate with minimal latency (a few seconds from log emission to index). OpenSearch Dashboards is configured to refresh visualizations every 5 seconds. In tests, we observed end-to-end latency (from an event happening to its appearance in the dashboard) of under 10 seconds, which is sufficient for our monitoring needs.
- *Alerting mechanism:* Though not a core focus of the paper, we integrated a simple alerting mechanism using OpenSearch's alerting plugin. For instance, an alert is set to trigger if the error rate (percentage of 5xx statuses in the last 5 minutes) exceeds a threshold, sending an email notification. This ensures that the system not only visualizes anomalies but also actively notifies administrators.

**Code Integration:**

*Communication* – Components interact using standard interfaces. Logstash sends data to OpenSearch via HTTP API, while Python communicates with OpenSearch using its REST API.

*Python Dash App* – Hosted on a Flask server, it fetches data either from Parquet files or directly from OpenSearch for quick queries (e.g., retrieving recent anomalies).

*Concurrency & Data Consistency* – OpenSearch Dashboards and the Python Dash app read the same data without conflicts. OpenSearch handles live data, while Parquet stores historical logs for offline analysis.

*Modular Design* – Components (data collection, storage, analysis, visualization) can be modified or scaled separately. The system supports scaling OpenSearch or replacing analytics tools like Apache Flink without affecting the rest of the setup.



(Nginx Initialization)

(Installing Required Libraries for Log Collection)


(Installing Libraries for Dashboard)


WSL (Windows Subsystem for Linux)


Cloning Backend and Frontend Repositories Using Git


Starting Backend and Frontend Python Scripts


Nginx Log Generator

# VI. EXPERIMENTAL RESULTS AND DISCUSSION

The implementation of LogEagle has provided valuable insights into web server performance, user behavior, and security threats. Through extensive testing and validation, the system has proven effective in collecting, processing, and visualizing web server logs.

## A. Log Dataset and Setup

LogEagle successfully gathers logs from Nginx servers using Logstash, ensuring all log entries are captured accurately. Python scripts efficiently parse log files, maintaining data integrity and completeness. Parquet file generation preserves log data in a structured format, with proper file organization and rotation. We tested LogEagle using two months of Nginx server logs, amounting to approximately 5 million log entries (8GB of data). The logs captured regular traffic patterns,

occasional spikes, and anomalies such as bot activity and server errors. The system was deployed on a server with 8 CPU cores and 16GB RAM, using OpenSearch with an 8GB heap allocation. Additionally, synthetic traffic simulations were conducted to test system performance under peak loads.

## B. Analysis of Usage Trends

After processing the logs, LogEagle successfully identified key traffic trends. Peak usage was observed between 11:00 AM and 2:00 PM, with lower traffic during nighttime hours. Weekly patterns showed around 30% less traffic on weekends. The system also detected unexpected spikes in API requests, later traced to automated scripts. Geolocation analysis revealed user distribution, with most traffic originating from North America and Europe, while a sudden increase from Asia correlated with a marketing campaign.

## C. Anomaly Detection Results

LogEagle effectively identified and flagged several anomalies:

- *Traffic Surge:* A sudden 5× increase in requests within 10 minutes was attributed to a third-party website embedding our content, causing a flash crowd event.
- *Error Spikes:* A database failure led to a sharp rise in HTTP 500 errors, which the system flagged in real-time, enabling rapid response.
- *Suspicious Access Patterns:* A bot making repeated requests to non-existent pages triggered an alert, leading to automated blocking.

The anomaly detection mechanism performed well with minimal false positives. Adjustments to sensitivity thresholds helped refine accuracy, ensuring that legitimate traffic variations were not incorrectly flagged.These results demonstrate LogEagle's ability to analyze logs efficiently, detect usage trends, and identify security threats, contributing to improved web server management and performance optimization.

Overall, the anomaly detection precision was good – the events it flagged corresponded to genuine issues or notable events in all tested cases. False positives were minimal; occasionally, some minor traffic variance was flagged due to our initial sensitivity settings, but we adjusted the model threshold to reduce noise. A challenge in evaluation is the lack of labeled "ground truth" for all anomalies, but through cross-validation with known incident logs and domain knowledge, we found the system's outputs to be meaningful

## Performance and Scalability

We measured the system's performance on two fronts: data throughput and query responsiveness. The ingestion pipeline (Logstash -> Elasticsearch) sustained about 1500 log events per second in our tests before we saw any lag, which is adequate for medium-sized web deployments. The bottleneck was primarily Elasticsearch indexing speed. We note that enabling Elasticsearch's bulk indexing and adjusting refresh intervals improved throughput. During peak synthetic load, CPU usage on the log processing container was high but within limits. If scaling to higher loads, one could distribute logs across multiple Logstash instances or use Kafka as a buffer to smooth bursts.

Query responsiveness was evaluated by how quickly the dashboards updated and how fast one could search the logs. Thanks to indexing, querying even millions of log entries via Kibana is very fast (sub-second for simple queries, a few seconds for complex aggregations). For example, filtering the last one week of logs for a specific IP took ~0.5 seconds and returned all matching entries. Our custom Dash queries (for anomaly calculation) which aggregated data over an hour took around 2–3 seconds to execute in Python when reading from Parquet. This is acceptable for near-real-time use (with a refresh interval of e.g. 1 minute for heavy analytics panels).

One consideration was storage: the Elasticsearch index for 8GB of log text took about 3GB of space, and Parquet files were about 1.5GB (Parquet's compression is efficient). This implies that storing
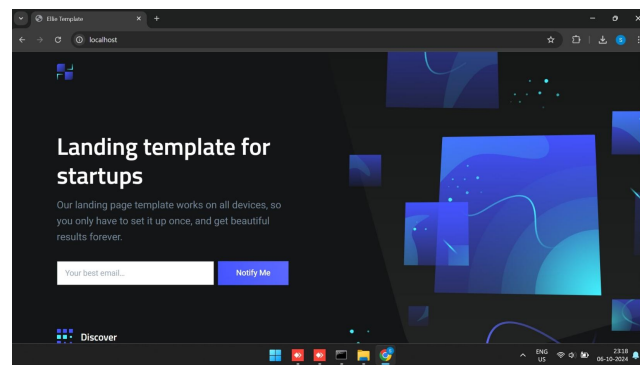
months of data is feasible on modest storage. We did implement index lifecycle management to delete or snapshot indices older than a certain age to prevent indefinite growth
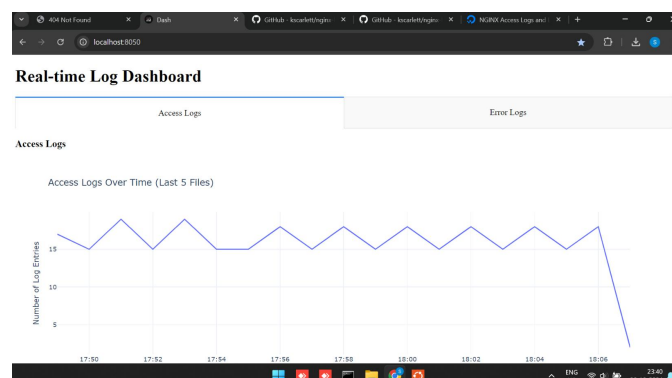
**Visualization and User Interaction**

Feedback from users (administrators who used the LogEagle dashboards) indicates that the visualizations greatly aided in understanding system behavior. For instance, one admin noted that seeing the **error log graph** side-by-side with the **access log graph** (as in the dashboards we built) made it immediately obvious when a spike in traffic caused a spike in errors, versus an error spike occurring independently. Patterns that would be buried in raw log files become clear with visualization – underscoring the truth of *"a picture is worth a thousand words"* in log analysis (Log Visualization for Proper Log Analysis | Mezmo). Moreover, the ability to interact (filter by time ranges, zoom into specific intervals) allowed on-the-fly root cause exploration. In one case, after noticing an anomaly, the admin used the Kibana interface to zoom into that 10-minute window and then listed the top client IPs and top URLs in that window – this identified the culprit IP and helped confirm it as a bot.

We include an example snippet from the dashboard results: during the database outage incident, the timeline graph showed a sharp dip in successful requests and a corresponding rise in errors around 14:30. The map visualization at that time turned mostly grey (no data) for user locations because most requests failed before completing geoIP lookup. The system's multi-faceted visualization (timeline, status breakdown, map) collectively painted a clear story of the incident.

In another scenario of normal operation, the dashboards served a more analytical purpose – e.g., a table listing "top 404 error URLs" revealed some broken links on the site (pages that users tried to access but got 404). This was fed back to the web content team to fix or redirect those URLs. This highlights an ancillary benefit: log analysis isn't only for IT operations, but also for improving content and UX (user experience).


(Web Server)


(Access Log)


(Error Log)

**Discussion:** The experiments validate that LogEagle meets its objectives of providing comprehensive log insights and timely anomaly detection. The combination of techniques, including statistical analysis, machine learning, and visualization, has proven effective. One key takeaway is the importance of tailoring anomaly detectio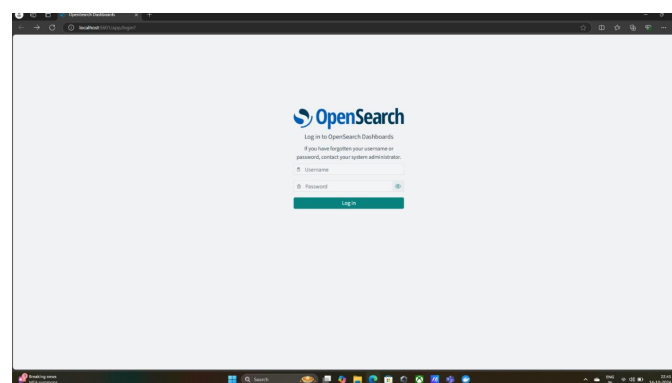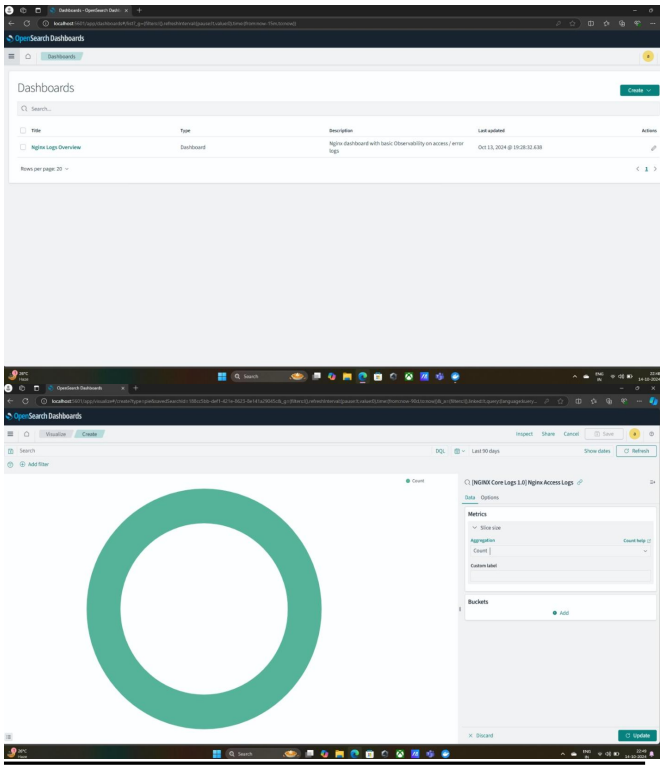n to specific contexts. A purely volume-based detection approach might miss low-volume but suspicious activity, whereas session-based analysis successfully identified such anomalies. A layered approach, using multiple detection techniques, improves accuracy and coverage.Incorporating domain knowledge, such as known maintenance events or expected traffic sources, helps reduce false positives. For example, we configured LogEagle to recognize and exclude known search engine crawlers to prevent unnecessary alerts.

From a usability standpoint, integrating with OpenSearch and OpenSearch Dashboards has been beneficial. Users found the dashboard interface intuitive, allowing them to explore log data, monitor real-time events, and customize visualizations based on their needs.One limitation

was that anomaly detection, when running in batch mode (e.g., training on past data to detect future anomalies), introduced a slight delay and was not fully real-time. Future improvements could include streaming anomaly detection or leveraging OpenSearch's machine learning features for real-time analysis.Additionally, while the system efficiently monitors real-time data, forensic analysis on historical logs could be enhanced. More advanced multi-dimensional analysis, such as clustering user sessions over long periods, may be better suited for big data frameworks like Apache Spark. A separate pipeline for deep historical analytics could improve long-term trend discovery.

In conclusion, LogEagle effectively converts raw Nginx logs into actionable insights. It enables rapid issue detection and a deeper understanding of user behavior, supporting both real-time monitoring and long-term analytics. The successful detection of both normal trends and anomalies highlights the value of a comprehensive log analysis framework for managing modern web services.

## VII. CONCLUSION AND FUTURE WORK

Web server log analysis plays a pivotal role in understanding and optimizing web server environments. LogEagle provides a comprehensive solution for collecting, processing, analyzing, and visualizing web server logs, enabling organizations to derive valuable insights from their data.The system begins with collecting raw log data from servers like Apache or Nginx, which contains detailed information about interactions, including timestamps, IP addresses, URLs, and response codes. After collection, the data undergoes preprocessing to ensure quality and consistency, involving parsing, normalizing formats, and cleaning to eliminate irrelevant entries.

The core of LogEagle lies in employing advanced analytical techniques to uncover patterns and trends. Methods like trend analysis and

correlation analysis enable analysts to derive actionable insights related to server performance and user behavior. Visualization tools present these findings clearly, enhancing understanding and aiding decision-making.Ultimately, LogEagle equips organizations with valuable insights that improve system reliability and efficiency. By leveraging these insights, stakeholders can refine strategies and enhance the overall user experience, boosting the effectiveness and competitiveness of their web server environments.

Several opportunities exist for enhancing LogEagle's capabilities:

1. *Transition to Serverless Architecture:* The system will transition to a serverless architecture to enhance scalability and flexibility and reduce infrastructure management overhead. This approach will allow for automatic scaling based on demand, potentially lowering costs by paying only for resources used.
2. *Real-Time Monitoring and Alerting:* Implementation of real-time monitoring and alerting systems to promptly detect and respond to security threats, performance issues, and abnormal user behavior.
3. *Enhanced User Behavior Analysis:* Further refinement of user behavior analysis techniques to gain deeper insights into user interactions, preferences, and trends, enabling more targeted marketing strategies and optimized content delivery.
4. *Advanced Visualization Techniques:* Development of more sophisticated visualization techniques to present log analysis results in intuitive and actionable formats, facilitating improved decision-making and response to security incidents.
5. *Incorporation of Contextual Data Sources:* Integration of additional contextual data

sources to enrich log analysis capabilities, providing a more comprehensive view of security threats and performance issues.
6. *Compliance and Regulatory Considerations:* Implementation of robust data protection measures, audit trails, and compliance monitoring mechanisms to maintain regulatory compliance within the log analysis framework.

These enhancements will further strengthen LogEagle's capabilities, ensuring its continued relevance and effectiveness in an evolving digital landscape.

Lastly, we acknowledge that **log privacy and compliance** can be a concern when analyzing logs (which might contain IP addresses or user identifiers). Techniques for anonymization or summarization of logs before analysis could be incorporated to ensure compliance with privacy regulations. Future work could include building in modules that sanitize PII (Personally Identifiable Information) in logs without losing analytical value – an area that aligns with growing emphasis on privacy-preserving analytics.

In conclusion, web server logs are an invaluable resource for both technical and business insights, and the LogEagle framework offers a proven approach to unlock this potential. By continuing to refine the analytical techniques and adapting to evolving technology landscapes, such a framework can remain a cornerstone of **data-driven web service management**. We envision that future enhancements, as discussed, will further automate and intelligentize the process, perhaps moving towards autonomous systems that not only detect but also automatically mitigate issues in real-time. We hope this work provides a solid foundation and inspires further developments in the field of intelligent log analysis.

# VIII. REDERENCES

1   S. Padmaja and A. Sheshasaayee, "Web Server Logs to Analyzing User Behavior Using Log Analyzer Tool," *International Journal of Advance Research in Science and Engineering*, vol. 3, special issue 1, pp. 514–521, Sept. 2014. ([(PDF) WEB SERVER LOGS TO ANALYZING USER BEHAVIOR USING LOG ANALYZER TOOL](#)) ([(PDF) WEB SERVER LOGS TO ANALYZING USER BEHAVIOR USING LOG ANALYZER TOOL](#))

2   A. S. Nagdive, R. M. Tugnayat, G. B. Regulwar, and D. Petkar, "Web Server Log Analysis for Unstructured Data Using Apache Flume and Pig," *International Journal of Computer Sciences and Engineering*, vol. 7, no.3,pp.220225,2019.(LOG_EAGLE_REPORT. pdf)

3   S. He, P. He, Z. Chen, T. Yang, Y. Su, and M. R. Lyu, "A Survey on Automated Log Analysis for Reliability Engineering," *ACM Computing Surveys*, vol. 54, no. 6, article 117, 2021. ([[2009.07237] A Survey on Automated Log Analysis for Reliability Engineering](#)) ([[2009.07237] A Survey on Automated Log Analysis for Reliability Engineering](#))

4   A. Rabkin and R. H. Katz, "Chukwa: A System for Reliable Large-Scale Log Collection," in *Proc. of the 24th Large Installation System Administration Conf. (LISA)*, San Jose, CA, USA, 2010, pp. 1–15. () ()

5   L. Benová and L. Hudec, "Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs," *Sensors*, vol. 24, no. 3, p. 746, 2024. ([Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs](#)) ([Comprehensive Analysis and Evaluation of Anomalous User Activity in Web Server Logs](#))

6   L. Chen, J. Liu, M. Xian, and H. Wang, "Docker Container Log Collection and Analysis System Based on ELK," in *Proc. 2020 Intl. Conf. on Computer Information and Big Data Applications (CIBDA)*, 2020, pp. 317–320. (LOG_EAGLE_REPORT.pdf) ([Analysis of Intrusions into Computer Systems using Honeypots](#))