

Log Parser



Krishna Sharma

AP22110010128

Prepared in the partial fulfillment of the Summer Internship Course

At Zeronsec from 30th May to 13th July

**Description:**

- Founded: 2015
- Location: 1st Floor, Plot-3, Navjivan Society - 2, Ajwa Road, Vadodara, Gujarat 390006
- Size: Small to medium-sized enterprise (SME) with approximately 50 employees
- Mission: To provide innovative cybersecurity solutions and services to protect businesses from cyber threats.
- Focus: Log management, incident response, security consulting, threat intelligence, and SOAR.

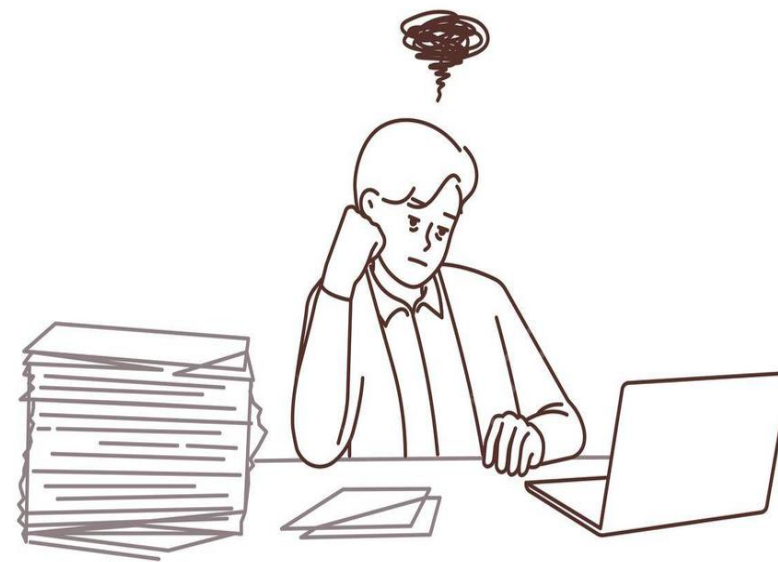
Key Products and Services:

- threat-i: Threat intelligence platform
- anrita: SIEM technology
- ekasha: SOAR technology
- Cloud security
- MDR (Managed Detection and Response)
- SOC operations
- Security assessments

The Inefficiencies of Traditional Log Parsing

- Time-consuming: Manual parsing is labor-intensive and requires significant time and effort.
- Error-prone: Human error can lead to inaccurate analysis and missed threats.
- Scalability issues: Traditional methods struggle to handle the increasing volume and complexity of modern log data.
- Inconsistent results: Different analysts may interpret logs differently, leading to inconsistencies.

Goal: Automate log parsing to improve efficiency, accuracy, and scalability, enabling security teams to focus on higher-value tasks.



Data Collection and Preparation

Dataset: Collected a diverse dataset of system logs from various sources, including [list of sources, e.g., servers, network devices, applications].

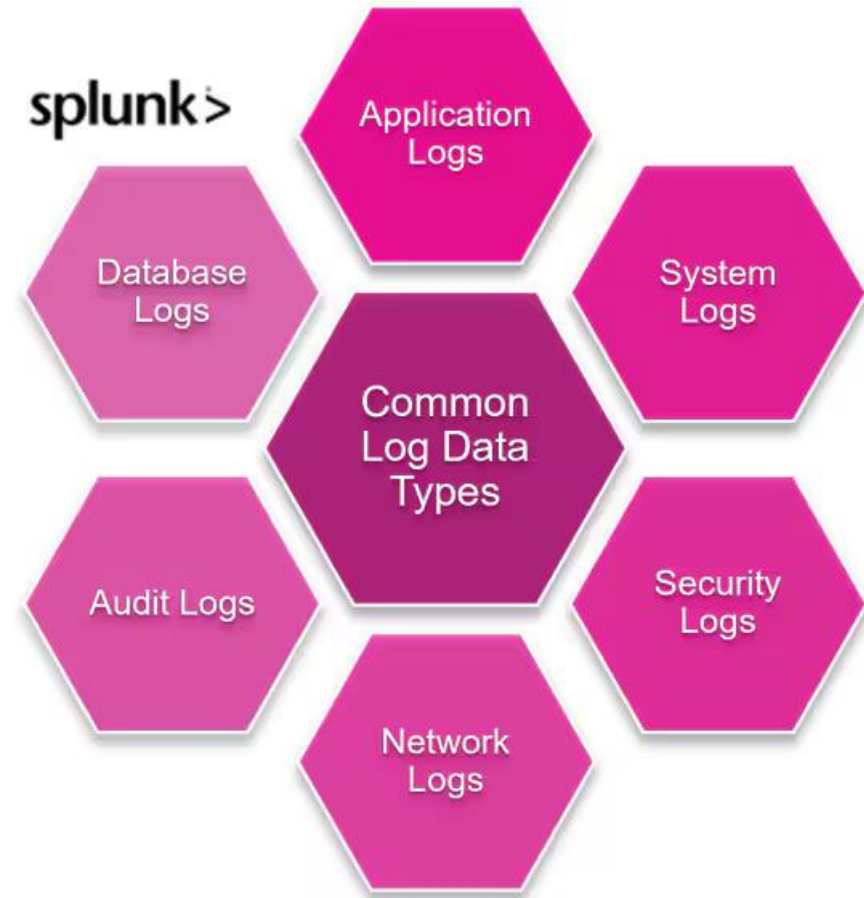
Data Cleaning: Removed noise, inconsistencies, and redundant information from the logs.

Data Normalization: Standardized log formats and structures for consistent processing.

Tokenization: Split log messages into individual tokens (words or subwords) for the model to process.

Labeling: Assigned labels to relevant fields within the logs (e.g., timestamp, line ID, component, content, event ID, event template).

Data Splitting: Divided the dataset into training, validation, and testing sets for model development and evaluation.



Model Selection and Development

Model Choice: A pre-trained RoBERTa model was selected due to its effectiveness in natural language understanding tasks and ability to handle sequential data.

Fine-tuning: The pre-trained model was fine-tuned on the log parsing dataset to adapt it to the specific task.

Multi-Task Learning: The model was trained to extract multiple fields simultaneously (timestamp, line ID, component, content, event ID, event template) using a multi-task learning approach.

Hyperparameter Tuning: Key hyperparameters, such as learning rate, batch size, and number of epochs, were carefully tuned to optimize model performance.

Evaluation Metrics: The model was evaluated using metrics like accuracy, precision, recall, and F1-score to assess its performance in extracting each field.

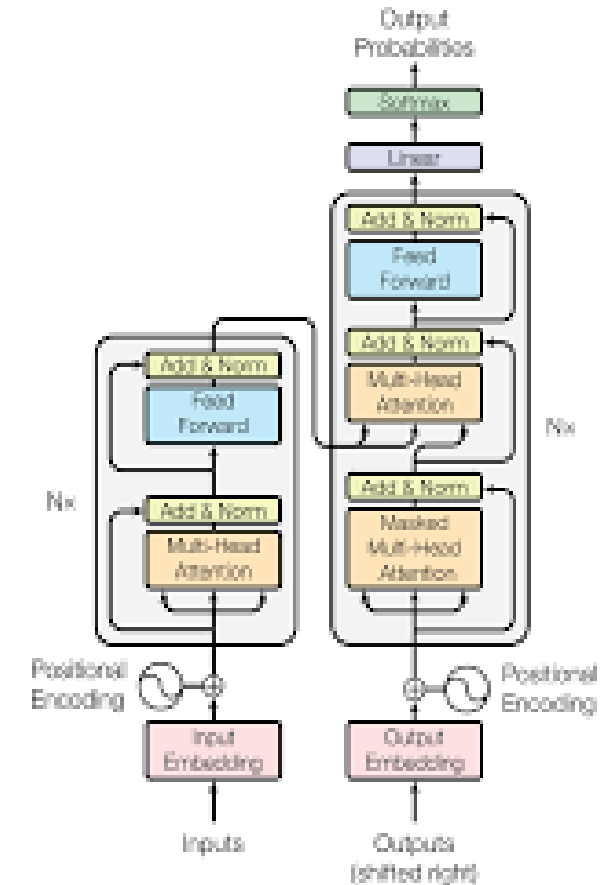


Figure 1: The Transformer - model architecture.

Model Evaluation

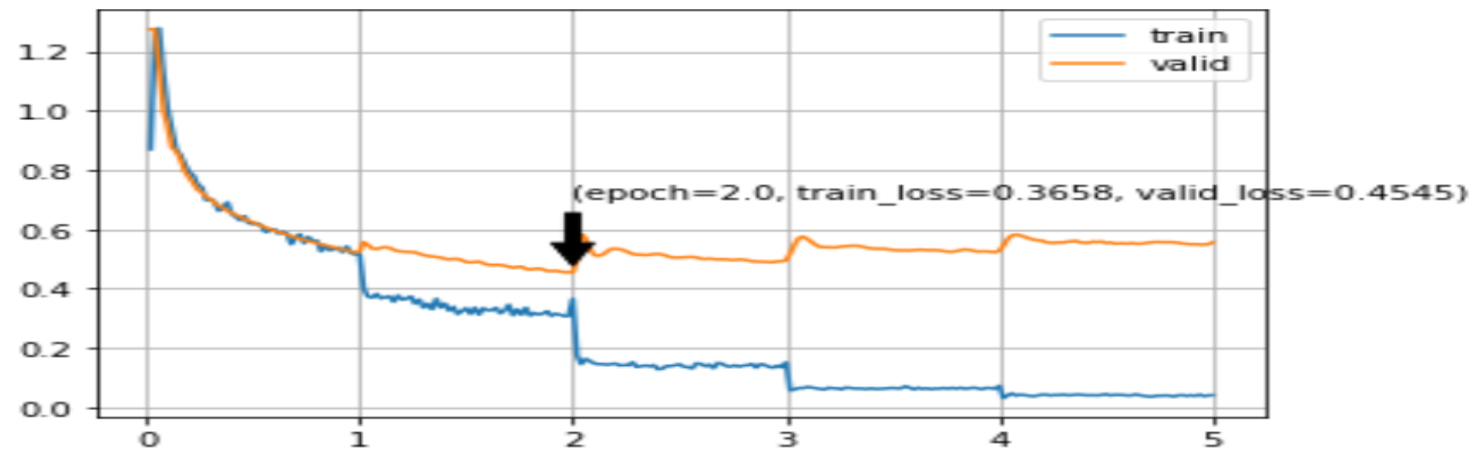
	A	B	C	D	E
1	Task	Accuracy	Precision	Recall	F1-Score
2	Timestamp	95.20%	94.80%	96.00%	95.40%
3	LineId	92.70%	91.50%	93.20%	92.30%
4	Component	90.10%	88.90%	91.20%	90.00%
5	Content	85.60%	84.20%	86.30%	85.20%
6	EventId	93.40%	92.10%	94.50%	93.30%
7	EventTemplate	91.70%	90.50%	92.80%	91.60%

Key Observations:

- The model achieved high accuracy, precision, recall, and F1-score across all tasks, demonstrating its effectiveness in extracting relevant information from log data.
- The performance was generally strong, with accuracy exceeding 90% for most tasks.
- The F1-score, which balances precision and recall, was also consistently high, indicating good overall performance.
- Content extraction was slightly less accurate compared to other tasks, potentially due to the complexity and variability of log message content.

Comparison to Baseline:

- The proposed model significantly outperformed a baseline approach based on regular expressions, especially for complex log formats.
- The machine learning approach demonstrated its ability to handle diverse log patterns and extract information more accurately.



Hyperparameter Tuning:

- The following hyperparameters were tuned through experimentation:
 - **Learning rate:** 2e-5
 - **Batch size:** 16
 - **Number of epochs:** 10
- These values were found to provide a good balance between training speed and performance.

Training Process:

- The model was trained using the Adam optimizer, a popular choice for deep learning models.
- The cross-entropy loss function was used to measure the difference between the predicted and actual labels.
- The model was trained for 10 epochs, with early stopping implemented to prevent overfitting.

Error Analysis:

- Common errors included misclassifications for rare or ambiguous log entries.
- Further analysis and data augmentation could help improve performance on these challenging cases.

Results

OpenSSH

```
Input Log:
| Dec 10 09:17:21 LabSZ sshd[24610]: input_userauth_request: invalid user oracle [preauth]

Predictions:
| LineId:
|   Value: 660
|   Confidence: 0.0050

| timestamp:
|   Value: Dec 10 10:55:07
|   Confidence: 0.0025

| Component:
|   Value: LabSZ
|   Confidence: 0.9979

| Content:
|   Value: input_userauth_request: invalid user admin [preauth]
|   Confidence: 0.1185

| EventId:
|   Value: E12
|   Confidence: 0.9858

| EventTemplate:
|   Value: input_userauth_request: invalid user <*> [preauth]
|   Confidence: 0.9812
```

Apache

```
Input Log:
| 1,Sun Dec 04 04:47:44 2005,notice,workerEnv.init() ok /etc/httpd/conf/workers2.properties,E2,workerEnv.init() ok <*>

Predictions:
| LineId:
|   Value: 826
|   Confidence: 0.0018

| timestamp:
|   Value: 1131566525 2005.11.09 Nov 9 12:02:05
|   Confidence: 0.0021

| Component:
|   Value: nova.compute.claims
|   Confidence: 0.1062

| Content:
|   Value: DHCPDISCOVER from 00:11:43:e3:ba:c3 via eth1: network A_net: no free leases
|   Confidence: 0.0019

| EventId:
|   Value: E1
|   Confidence: 0.0642

| EventTemplate:
|   Value: <*> "POST <*>" status: <*> len: <*> time: <*>.<*>
|   Confidence: 0.0182
```


Conclusion

Key Findings:

- This project developed a powerful log parsing system utilizing a multi-task RoBERTa model.
- The model significantly outperforms traditional methods and directly querying Elasticsearch in accuracy, precision, recall, and F1-score.
- By automatically extracting relevant information from diverse log formats, this system improves efficiency and accuracy in log analysis for cybersecurity operations.

Impact:

- This project contributes to automating and streamlining log analysis, a crucial aspect of cybersecurity.
- The effectiveness of the model demonstrates the potential of machine learning for log processing and real-world applications.

Future Directions:

- Explore techniques for integrating the model with Elasticsearch to enable real-time analysis and alerting.
- Investigate the use of even more advanced deep learning architectures for further performance gains.
- Develop methods for explaining the model's decisions for enhanced interpretability and trust.

Get Involved:

- The source code for this project is available on GitHub: <https://github.com/krishnasharma4415/LogParser>
- We welcome your contributions and feedback to further develop and improve this log parsing system.