# Log Summarisation for Defect Evolution Analysis

Rares Dolga
JPMorgan Chase
UK
rares.dolga@jpmchase.com

Ran Zmigrod
JPMorgan Chase
UK
ran.zmigrod@jpmorgan.com

Rui Silva
JPMorgan Chase
UK
rui.silva@jpmorgan.com

Salwa Alamir
JPMorgan Chase
UK
salwa.alamir@jpmchase.com

Sameena Shah
JPMorgan Chase
UK
sameena.shah@jpmchase.com

## ABSTRACT

Log analysis and monitoring are essential aspects in software maintenance and identifying defects. In particular, the temporal nature and vast size of log data leads to an interesting and important research question: How can logs be summarised and monitored over time? While this has been a fundamental topic of research in the software engineering community, work has typically focused on heuristic-, syntax-, or static-based methods. In this work, we suggest an online semantic-based clustering approach to error logs that dynamically updates the log clusters to enable monitoring code error life-cycles. We also introduce a novel metric to evaluate the performance of temporal log clusters. We test our system and evaluation metric with an industrial dataset and find that our solution outperforms similar systems. We hope that our work encourages further temporal exploration in defect datasets.

## CCS CONCEPTS

• **Theory of computation → Unsupervised learning and clustering**; • **Software and its engineering → Software defect analysis**.

## KEYWORDS

Defect Detection, Log Analysis, Clustering, NLP

## 1 INTRODUCTION

Log data is crucial in the analysis, maintenance, and fixing of errors in software systems [1, 4]. Logs are verbose documents that detail a system's state throughout its execution lifecycle and as

such, lead to vast amounts of data that make manual analysis intractable in practice [18]. Consequently, both research and industry communities have focused on automatic log analysis tasks such as failure diagnosis and prediction [32], anomaly detection [7], log comprehension [17], *inter alia*.

Log data for a system is temporal, i.e., logs stream sequentially as a system is developed and tested. Therefore, one can monitor the existence and persistence of code defects by monitoring their existence in logs [22]. To do this, past research has focused on focused log comprehension [19, 25] and log evolution tracking. The latter task has mainly been approached with heuristic-based systems [15, 28] and unsupervised syntax-based systems [10–12, 20]. These systems work well for small domains where data is structured, however, the unstructured and large nature of modern log data limit the success of these techniques for real world data.

Another approach to capture log evolution is to consider the semantic representations of logs. Using semantic rather than syntactic meaning for logs makes sense as they may vary across applications and developers (i.e., have different free-form text). Semantic representations for logs have been previously used for anomaly detection and defect prediction due to the free-text nature of logs [23, 25]. Moreover, past work has also focused on static semantic log clustering using common approaches (e.g., K-Means or Gaussian Mixture Models) [9]. While dynamic clustering techniques exist [2], they lead to overly erratic log evolution. Therefore, a gap exists in the software defect tracking literature that we aim to explore.

In this paper, we introduce a novel online algorithm for clustering and monitoring logs based on their semantic representation. Through collaboration with a team of 61 Software Reliability Engineers (SREs), who are responsible for almost 500 applications that jointly produce roughly 100,000 logs per day, we extracted three key criteria for successful monitoring of log evolution: identifying errors, providing meaningful representative messages, and minimising disruption of log clusters. These criteria were used to construct a novel evaluation metric that captures real world performance for temporal log clustering. Experiments using a private industrial dataset demonstrate that our online clustering method outperforms existing static and dynamic clustering methods. Additionally, we demonstrate that using a richer semantic representation further improves performance. We note that our work could be used in generating more high level error log topics in both temporal and non-temporal datasets. The contribution of this works are:

(1) We introduce a novel online semantic-based clustering algorithm for classifying errors and monitoring defect evolution.

(2) We construct a performance metric for log evolution based on the experience of real SREs.

(3) We demonstrate that our algorithm improves upon past clustering techniques on a private dataset and provide baseline metrics on public log datasets.

## 2 RELATED WORK

Log mining in the context of software defects is a well-researched topic [19] which mainly focuses on failure diagnosis [32], anomaly detection [7] and log comprehension [17]. Older research focuses on tracking log patterns over time through heuristics like frequency, sentence length or n-grams to group logs together [3, 15, 28]. More recently, the FLAP system was introduced as an end-to-end framework for log analysis [22] based on syntactic log properties and an iterative clustering process [24]. Other work exists that utilises syntactic features to group logs [8, 11, 29, 30] and track formed clusters [5, 12]. The techniques presented in these works all fail to extract a representative log for clustering which is a desirable quality in terms of the explainability of the log error types. Similar work has shown that it is possible to achieve syntactic representatives [10, 14, 20]. Furthermore, past work has also utilised genetic algorithms to extract log templates [26].

Semantic representations of logs have been explored in anomaly detection and defect prediction [6, 23, 25, 33]. Research also exists that examines log comprehension and clustering using semantic representations (e.g., Word2Vec) [9]. This work differs from the system proposed in this paper in two fundamental ways. Firstly, it treats logs using a bag-of-words (BoW) approach and so does not consider word order nor frequency. Therefore, it discards important information that exists in more complex representations such as the sentence embeddings of SBERT [27]. Secondly, it applies the DBSCAN clustering algorithm which is a static algorithm and so it is impossible to track log clusters over time. There are streaming clustering techniques such as DenStream [2] which work with vector representations, however, they do not provide representative extraction and can have an erratic convergence.

## 3 TRACKING ERROR LOG EVOLUTION

This paper concerns the monitoring of error logs which are descriptive of various defects encountered in a system's execution. We consider a **log** to be an individual message sent by the system during its execution (logs sent from the same execution are related by an execution ID). In order to motivate and devise an effective solution to our task, we introduce three key criteria for successful monitoring of log evolution that were reached in collaboration with over 60 SREs. We further motivate our criteria by a visual example in Figure 1.

**Identifying Errors over Time.** We wish to correctly discover patterns that represent a common defect or error expressed by logs. We expect the types of errors to change over time and be correctly captured by the algorithm. As can be seen in Figure 1, over time, we no longer witness "User Interface Error" defects while we begin to witness "Connection Timeout" defects.

**Capturing Meaningful Representatives.** It is important that we are able to describe each cluster by its common defect. As such, a

successful algorithm should provide a meaningful representative to each cluster that an end-user is able to understand. A meaningful representative is visually depicted to be near the centre of each cluster in Figure 1. Additionally, each cluster has an understandable and succinct representative error that aligns with common code defects as seen in the tables of Figure 1.

**Minimising Disruption.** In order to effectively monitor the life-cycle of an error or defect, clusters must remain active until their associated errors are fully resolved. As such, we require clusters to evolve smoothly over time while maintaining high quality clusters. There is minimal disruption illustrated in Figure 1 as cluster 1 expands and eventually breaks off into cluster 4 (perhaps along with cluster 2). The smooth evolution over time is needed to not introduce or remove a defect cluster too quickly or erratically.

### 3.1 A Novel Online Clustering Algorithm

One of the contributions of this work is a novel online clustering algorithm for error logs that is motivated by the above criteria. We draw some inspiration from DenStream [2] who create dynamically changing clusters. In particular, we adopt the same distance metric between logs and cluster centroids. The intuition of our online clustering algorithm is straight-forward: When new data points are received, we treat points either one by one or in batches and merge them into the cluster with the smallest cosine distance to the centroid of that cluster. We formulate our algorithm in Alg. 1 using three hyperparameters:

$\theta$: The acceptance threshold of a log to a cluster. If the distance between a log and all cluster centroids is less than $\theta$, a new cluster is created.

$\alpha$: The rate of cluster centroid evolution. We update cluster centroids as a rolling average based on new logs that are merged into a cluster. Consequently, cluster centroids naturally change over time. To prevent this from occurring too quickly (and thus minimise disruption), we use $\alpha$ to slow change.

$\gamma$: The minimum cluster size before a centroid of a cluster may change. We use this hyperparameter to ensure outliers do not heavily impact smaller clusters that are more liable to centroid changes.

For an end user to clearly interpret the clustering results, we extract an error message from each cluster. For each cluster, we use cosine similarity to select the sample which is closest to the centroid and use its log as the cluster's representative defect. Alternatively, the Levenshtein score can be used to get the average message [21]. While this method may give a better message as it is more specific to the log text, it is more computationally expensive and so we stick to the cosine similarity approach.[1]

### 3.2 Performance Metrics

We next introduce a novel metric that formalises the three key criteria identified at the start of this section. Our metric is devised of three scores: Silhouette score ($S$), Representative Similarity ($R$) and

---

[1]Let each cluster have $m$ logs, the average log length is $n$, and the semantic representation dimensionality is $d$ where $d << n$. Computing the Levenshtein score has a time complexity of $O(m^2 n^2)$ as opposed to $O(md)$ for the cosine similarity approach.
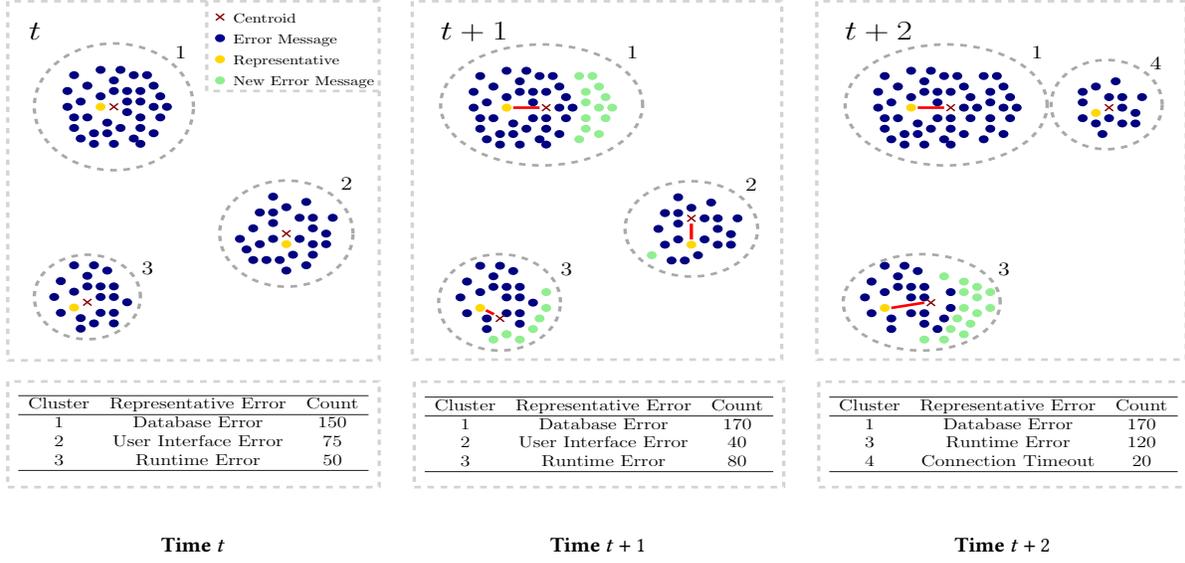
**Figure 1: Example of Error Log Evolution. The top row provides a visual aid of the clusters while the bottom row details the specific defects found at each point in time.**

**Data:** List of batches
**Hyperparameters:** $\theta$; $\alpha$; $\gamma$
clusters = [];
**for** *batch* **in** *data* **do**
    **for** *p* **in** *batch* **do**
        min_dist = inf; c = null;
        **for** *clust* **in** *clusters* **do**
            dist = 1 - cosine_sim(clust.cen, p);
            **if** *dist < min_dist* **then**
                min_dist = dist;
                c = clust;
            **end**
        **end**
        **if** *min_dist ≤ θ* **then**
            **if** *c.len ≥ γ* **then**
                c.cen = (1-$\alpha$) · c.cen + $\alpha p$;
            **else**
                c.cen = $\frac{c.len}{c.len+1}$ · c.cen + $\frac{1}{c.len+1}p$;
            **end**
            c.len += 1
        **else**
            clust = new Clust(cen=p, len=1);
            clusters.add(clust);
        **end**
    **end**
**end**

**Algorithm 1:** Our online clustering algorithm, with updates coming as individual points. Batched updates are more efficient.

Number of Clusters Similarity (*C*). Each of these scores is computed over batches and is in the range [0, 1].

The Silhouette score is a standard metric used for assessing unlabelled clustering. It provides a measure of the similarity of each log to its own cluster versus other clusters. Traditionally, the score is given in the range $[-1, 1]$, but we scale it so that it is comparable with our other scores.

$$S = \frac{1}{B} \sum_{b=1}^{B} \frac{Silhouette_b + 1}{2} \qquad (1)$$

where $Silhouette_b$ is the Silhouette score for bath $b$.

The representative similarity score compares the representatives of each cluster $c$ for batches $b$ and $b + 1$ coming in at times $t$ and $t + 1$ respectively. We use this score to evaluate the evolution of the representative logs for each cluster. We use the cosine similarity as our distance metric.

$$R = \frac{1}{B} \sum_{b=1}^{B} cosine\_sim(repres_b, repres_{b+1}) \qquad (2)$$

where $repres_b$ is the semantic representations for batch $b$.

Lastly, the number of clusters similarity score measures how smoothly the number of clusters evolves over time.

$$C = 1 - \frac{1}{B} \sum_{b=1}^{B} \frac{|nr\_clust_b - nr\_clust_{b-1}|}{max(nr\_clust_b, nr\_clust_{b-1})} \qquad (3)$$

where $nr\_clust_b$ is the number of clusters in batch $b$.

Our log cluster evolution metric (*LCE*) is then comprised of a linear combination of the three scores above, where all scores are scaled to be in the range [0, 1] and $w_S + w_R + w_C = 1$

$$LCE = w_S S + w_R R + w_C C \qquad (4)$$

## 4 EXPERIMENTS

In this section, we measure the performance of our online clustering algorithm using different semantic representations.

### 4.1 Data

We primarily conduct experiments on a private industrial dataset from the SREs we worked with. To do this, we collected log data from a database containing all monitored applications. We then eliminated variability by removing timestamps and URLs. The monitored applications generate around 100,000 logs per day, of which, about 1.5% are related to errors and defects. We extracted two months' worth of data, resulting in a total of 57,000 error logs. Each log is comprised of structured fields (e.g., system id, date, log level) as well as unstructured, free-form text.

Since our dataset is private and cannot be shared, we also evaluate our method on similar public datasets. We use Loghub [13], a collection of logs from different systems which have a date element, a log level and free-form text representing the meaning of the log. We note that this dataset is significantly simpler than the private industry dataset due to generally shorter and less varied log texts. We use a sample of 2,000 logs for each of the Loghub systems: HDFS_2, Linux, Zookeeper and OpenStack. Note that each dataset is much smaller than the private dataset and examines a single framework rather than many applications. Throughout both the public and private datasets, we consider a log to be stale after one month since its creation.

**Pre-processing and Semantic Representations.** Logs are a challenging data type for text understanding models because they tend to contain incomplete words, words with special characters (e.g., underscore), and system-specific terms. Consequently, we cleaned the data using standard natural languages processing techniques, such as tokenization[31], lemmatization[16], and removal of English stop words given by the SpaCy library[2]. Given the cleaned logs, we constructed semantic representations of the logs using average Word2Vec[3] vectors as well as SBERT sentence embeddings [27].

**Temporal Splits.** For the private dataset, we first split the two months' data into a one-month batch and the rest into five-day batches. The second split considers only one-day batches. The reason for these two separate temporal settings is to show results on stationary versus dynamic data. In the first split, due to the high amount of data in the snapshot, the underlying distribution does not differ from the next timestamps. However, for one-day splits, one batch is not a meaningful sample of the entire distribution of logs. We use these datasets to compare our model with other clustering algorithms. We use one day batches for the public datasets.

### 4.2 Models and Algorithms

To benchmark our algorithm against the literature we picked two well-known models which work with vector representations. Gaussian Mixture Model (GMM) is an offline clustering algorithm which was selected because it allows clusters of any shape to be formed. We adapt the algorithm to be online by re-initialising each batch

[2]https://spacy.io/
[3]https://spacy.io/usage/spacy-101/#vectors-similarity

**Table 1: Clustering evolution performance on private dataset.**

| Model | # Clusters | *R* | *S* | *C* | *LCE* |
|---|---|---|---|---|---|
| | 1 Month snapshot + 5 day batches | | | | |
| GMM$_{\text{Word2Vec}}$ | 11 | 0.964 | 0.768 | **1** | 0.9 |
| DenStream$_{\text{Word2Vec}}$ | 48 | 0.92 | 0.771 | 0.803 | 0.83 |
| Ours$_{\text{Word2Vec}}$ | 198 | 0.999 | 0.865 | 0.96 | 0.94 |
| Ours$_{\text{SBERT}}$ | 27 | **0.999** | **0.97** | **0.96** | **0.97** |
| | 1 day batches | | | | |
| GMM$_{\text{Word2Vec}}$ | 11 | 0.914 | 0.727 | 1 | 0.88 |
| DenStream$_{\text{Word2Vec}}$ | 48 | 0.92 | 0.771 | 0.862 | 0.85 |
| Ours$_{\text{Word2Vec}}$ | 205 | **0.999** | 0.879 | **0.999** | 0.95 |
| Ours$_{\text{SBERT}}$ | 28 | **0.999** | **0.972** | **0.998** | **0.98** |

**Table 2: Clustering evolution performance on private and public datasets using Ours$_{\text{SBERT}}$ and one day batches.**

| Dataset | # Clusters | *R* | *S* | *C* | *LCE* |
|---|---|---|---|---|---|
| Private Dataset | 28 | 0.999 | 0.972 | 0.998 | 0.980 |
| HDFS_2 | 18 | 0.990 | 0.990 | 0.990 | 0.990 |
| Linux | 31 | 0.990 | 0.860 | 0.980 | 0.940 |
| Zookeper | 19 | 0.990 | 0.990 | 0.990 | 0.990 |
| OpenStack | 14 | 0.990 | 0.790 | 0.990 | 0.920 |

with the mean and covariance of the model trained on the previous batch. We choose DenStream [2] as the second algorithm because it trains in an online fashion and also allows clusters of varying shapes. The optimal values for the hyper-parameters of all models were selected based on a grid search conducted over past data.[4] For each algorithm we only use Word2Vec (embedding dimension of 300) to compare the algorithm's performance. We further use SBERT (embedding dimension of 756) for our algorithm to show the improvement in using a more complex and rich semantic representation. Non eof the models were fine-tuned on the data.

### 4.3 Results

For all results, *LCE* was computed with equal contribution weights. Table 1 shows the performance results across the different algorithms and semantic representations on the private dataset. We observe that for both temporal settings, our approach matches or outperforms both the GMM and DenStream algorithms when using Word2Vec semantic representations across each individual score as well as *LCE*. We note that GMM, the static clustering algorithm, performs better when it has the initial one month snapshot whereas the dynamic approaches perform consistently across both settings. This is expected as GMM requires the initial mean and variance to attain better performance.

The biggest improvement between our approach DenStream is observed in the cluster number similarity score. DenStream exhibits an erratic evolution which we are able to mitigate in our

[4]For our algorithm, the hyperparameters were $\theta = 0.05$, $\alpha = 0.1$, and $\gamma = 100$.

approach using our evolution rate ($\alpha$) and minimum cluster size hyperparameters ($\gamma$). The strong and consistent performance of our approach across both temporal settings suggests that our algorithm adapts easily to both stationary and non-stationary streams. The robustness of our algorithm is further shown in Table 2 where we achieve good results across public log datasets.[5]

## 5 THREATS TO VALIDITY

We note a few threats to validity in this work. Firstly, we must be cautious of overfitting to our data which would retract from the generalisibility of the model. As clustering is an unsupervised problem (i.e., no labelled data), we cannot use metrics such as ARI to evaluate our predictions. While the Silhouette score and manual qualitative analysis can be done to demonstrated good clusters, having access to a small labelled dataset would allow for more robust evaluation. This is currently a gap in log-based software defect datasets that we hope is addressed in future work. Indeed, one could use a model as presented in this work to generate data.

## 6 CONCLUSION

In this paper, we introduced a novel online semantic-based clustering algorithm for error logs. Our algorithm is able to cluster log message streams, track the reduction and emergence of defects that arise in code, and provide a representative defect message for each cluster. The algorithm has configurable cluster granularity and cluster evolution rate through hyperparameters. Furthermore, we introduced a novel metric named *LCE* that can evaluate defect tracking performance based on industry specified criteria. We demonstrated that our algorithm outperforms existing log clustering approaches using a private dataset and showed strong performance across public datasets. We hope that this work encourages more work in defect detection to examine the temporal properties of logs and incorporate it into modern systems.

## DISCLAIMER

This paper was prepared for informational purposes by the Artificial Intelligence Research group of JPMorgan Chase & Co and its affiliates ("JP Morgan"), and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

## REFERENCES

[1] Betsy Beyer, Chris Jones, Jennifer Petoff, and Niall Richard Murphy. 2016. *Site Reliability Engineering: How Google Runs Production Systems.* http://landing.google.com/sre/book.html

[2] Feng Cao, Martin Ester, Weining Qian, and Aoying Zhou. 2006. Density-Based Clustering over an Evolving Data Stream with Noise. In *SDM*.

[3] Hetong Dai, Heng Li, Che-Shao Chen, Weiyi Shang, and Tse-Hsun Chen. 2022. Logram: Efficient Log Parsing Using $n$-Gram Dictionaries. *IEEE Transactions on Software Engineering* 48, 3 (2022), 879–892. https://doi.org/10.1109/TSE.2020.3007554

[4] Steven Davies and Marc Roper. 2013. Bug localisation through diverse sources of information. In *2013 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW)*. 126–131. https://doi.org/10.1109/ISSREW.2013.6688891

[5] Min Du and Feifei Li. 2019. Spell: Online Streaming Parsing of Large Unstructured System Logs. *IEEE Transactions on Knowledge and Data Engineering* 31, 11 (2019), 2213–2227. https://doi.org/10.1109/TKDE.2018.2875442

[6] Min Du, Feifei Li, Guineng Zheng, and Vivek Srikumar. 2017. DeepLog: Anomaly Detection and Diagnosis from System Logs through Deep Learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security* (Dallas, Texas, USA) *(CCS '17)*. Association for Computing Machinery, New York, NY, USA, 1285–1298. https://doi.org/10.1145/3133956.3134015

[7] Qiang Fu, Jian-Guang Lou, Yi Wang, and Jiang Li. 2009. Execution Anomaly Detection in Distributed Systems through Unstructured Log Analysis. In *2009 Ninth IEEE International Conference on Data Mining*. 149–158. https://doi.org/10.1109/ICDM.2009.60

[8] Qiang Fu, Jian-Guang Lou, Yi Wang, and Jiang Li. 2009. Execution Anomaly Detection in Distributed Systems through Unstructured Log Analysis. *2009 Ninth IEEE International Conference on Data Mining* (2009), 149–158.

[9] Maria Grigorieva and Dmitry Grin. 2021. Clustering error messages produced by distributed computing infrastructure during the processing of high energy physics data. *International Journal of Modern Physics A* 36 (04 2021), 2150070. https://doi.org/10.1142/S0217751X21500706

[10] Nentawe Gurumdimma, Arshad Jhumka, Maria Liakata, Edward Chuah, and James Browne. 2015. Towards Detecting Patterns in Failure Logs of Large-Scale Distributed Systems. In *2015 IEEE International Parallel and Distributed Processing Symposium Workshop*. 1052–1061. https://doi.org/10.1109/IPDPSW.2015.109

[11] Hossein Hamooni, Biplob Debnath, Jianwu Xu, Hui Zhang, Guofei Jiang, and Abdullah Mueen. 2016. LogMine: Fast Pattern Recognition for Log Analytics. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management* (Indianapolis, Indiana, USA) *(CIKM '16)*. Association for Computing Machinery, New York, NY, USA, 1573–1582. https://doi.org/10.1145/2983323.2983358

[12] Pinjia He, Jieming Zhu, Zibin Zheng, and Michael R. Lyu. 2017. Drain: An Online Log Parsing Approach with Fixed Depth Tree. In *2017 IEEE International Conference on Web Services (ICWS)*. 33–40. https://doi.org/10.1109/ICWS.2017.13

[13] Shilin He, Jieming Zhu, Pinjia He, and Michael R. Lyu. 2020. Loghub: A Large Collection of System Log Datasets towards Automated Log Analytics. https://doi.org/10.48550/ARXIV.2008.06448

[14] P.W.D.C. Jayathilake, N. R. Weeraddana, and H. K. E. P. Hettiarachchi. 2017. Automatic detection of multi-line templates in software log files. In *2017 Seventeenth International Conference on Advances in ICT for Emerging Regions (ICTer)*. 1–8. https://doi.org/10.1109/ICTER.2017.8257824

[15] Yexi Jiang, Chang-Shing Perng, and Tao Li. 2011. Natural Event Summarization. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management* (Glasgow, Scotland, UK) *(CIKM '11)*. Association for Computing Machinery, New York, NY, USA, 765–774. https://doi.org/10.1145/2063576.2063688

[16] Divya Khyani and Siddhartha B S. 2021. An Interpretation of Lemmatization and Stemming in Natural Language Processing. *Shanghai Ligong Daxue Xuebao/Journal of University of Shanghai for Science and Technology* 22 (01 2021), 350–357.

[17] Jerry Kiernan and Evimaria Terzi. 2009. Constructing Comprehensive Summaries of Large Event Sequences. *ACM Trans. Knowl. Discov. Data* 3, 4, Article 21 (dec 2009), 31 pages. https://doi.org/10.1145/1631162.1631169

[18] Max Landauer, Florian Skopik, Markus Wurzenberger, and Andreas Rauber. 2020. System log clustering approaches for cyber security applications: A survey. *Computers Security* 92 (2020), 101739. https://doi.org/10.1016/j.cose.2020.101739

[19] Max Landauer, Florian Skopik, Markus Wurzenberger, and Andreas Rauber. 2020. System log clustering approaches for cyber security applications: A survey. *Computers and Security* 92 (2020), 101739. https://doi.org/10.1016/j.cose.2020.101739

[20] Laetitia Leichtnam, Eric Totel, Nicolas Prigent, and Ludovic Me. 2017. STARLORD: Linked security data exploration in a 3D graph. In *2017 IEEE Symposium on Visualization for Cyber Security (VizSec)*. 1–4. https://doi.org/10.1109/VIZSEC.2017.8062203

[21] Vladimir I. Levenshtein. 1965. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet physics. Doklady* 10 (1965), 707–710.

[22] Tao Li, Yexi Jiang, Chunqiu Zeng, Bin Xia, Zheng Liu, Wubai Zhou, Xiaolong Zhu, Wentao Wang, Liang Zhang, Jun Wu, Li Xue, and Dewei Bao. 2017. FLAP: An End-to-End Event Log Analysis Platform for System Management. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data*

---

[5]We further conducted a brief qualitative analysis with the SREs we worked with. They provided positive feedback with regards to the cluster evolution over time and pointed out certain clusters matched specific errors they were aware of.

*Mining* (Halifax, NS, Canada) *(KDD '17)*. Association for Computing Machinery, New York, NY, USA, 1547–1556. https://doi.org/10.1145/3097983.3098022

[23] Xiaoyun Li, Pengfei Chen, Linxiao Jing, Zilong He, and Guangba Yu. 2020. Swiss-Log: Robust and Unified Deep Learning Based Log Anomaly Detection for Diverse Faults. In *2020 IEEE 31st International Symposium on Software Reliability Engineering (ISSRE)*. 92–103. https://doi.org/10.1109/ISSRE5003.2020.00018

[24] Adetokunbo Makanju, Ayse Nur Zincir-Heywood, and Evangelos E. Milios. 2009. Clustering event logs using iterative partitioning. In *Knowledge Discovery and Data Mining*.

[25] Weibin Meng, Ying Liu, Yichen Zhu, Shenglin Zhang, Dan Pei, Yuqing Liu, Yihao Chen, Ruizhi Zhang, Shimin Tao, Pei Sun, and Rong Zhou. 2019. LogAnomaly: Unsupervised Detection of Sequential and Quantitative Anomalies in Unstructured Logs. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, 4739–4745. https://doi.org/10.24963/ijcai.2019/658

[26] Salma Messaoudi, Annibale Panichella, Domenico Bianculli, Lionel Briand, and Raimondas Sasnauskas. 2018. A Search-Based Approach for Accurate Identification of Log Message Formats. In *2018 IEEE/ACM 26th International Conference on Program Comprehension (ICPC)*. 167–16710.

[27] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. https://doi.org/10.48550/ARXIV.1908.10084

[28] Keiichi Shima. 2016. Length Matters: Clustering System Log Messages using Length of Words. https://doi.org/10.48550/ARXIV.1611.03213

[29] Liang Tang, Tao Li, and Chang-Shing Perng. 2011. LogSig: Generating System Events from Raw Textual Logs. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management* (Glasgow, Scotland, UK) *(CIKM '11)*. Association for Computing Machinery, New York, NY, USA, 785–794. https://doi.org/10.1145/2063576.2063690

[30] Risto Vaarandi and Mauno Pihelgas. 2015. LogCluster - A data clustering and pattern mining algorithm for event logs. In *2015 11th International Conference on Network and Service Management (CNSM)*. 1–7. https://doi.org/10.1109/CNSM.2015.7367331

[31] Jonathan Webster and Chunyu Kit. 1992. Tokenization as the initial phase in NLP. 1106–1110. https://doi.org/10.3115/992424.992434

[32] Shenglin Zhang, Ying Liu, Weibin Meng, Zhiling Luo, Jiahao Bu, Sen Yang, Peixian Liang, Dan Pei, Jun Xu, Yuzhi Zhang, Yu Chen, Hui Dong, Xianping Qu, and Lei Song. 2018. PreFix: Switch Failure Prediction in Datacenter Networks. *Abstracts of the 2018 ACM International Conference on Measurement and Modeling of Computer Systems* (2018).

[33] Xu Zhang, Yong Xu, Qingwei Lin, Bo Qiao, Hongyu Zhang, Yingnong Dang, Chunyu Xie, Xinsheng Yang, Qian Cheng, Ze Li, Junjie Chen, Xiaoting He, Randolph Yao, Jian-Guang Lou, Murali Chintalapati, Furao Shen, and Dongmei Zhang. 2019. Robust Log-Based Anomaly Detection on Unstable Log Data. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (Tallinn, Estonia) *(ESEC/FSE 2019)*. Association for Computing Machinery, New York, NY, USA, 807–817. https://doi.org/10.1145/3338906.3338931