

# Malware Search Engine

Ankit Patil  
School of Computer Science & IT  
Symbiosis Skills & Professional  
University  
Pune, India  
ankitpatil925@gmail.com

Darshan Rawal  
School of Computer Science & IT  
Symbiosis Skills & Professional  
University  
Pune, India  
drawal1713@gmail.com

Adnan Khan  
School of Computer Science & IT  
Symbiosis Skills & Professional  
University  
Pune, India  
adnanskhan56@gmail.com

Guide

Prof. Rahul Raut  
School of Computer Science & IT  
Symbiosis Skills & Professional  
University  
Pune, India  
rahul.raut@sspu.ac.in

## *Abstract—*

The spread of malware poses a major cybersecurity threat and requires effective detection, analysis and classification. This article introduces a Malware Search Engine (MSE) that solves these problems. Leveraging advanced information retrieval, machine learning, and data analysis techniques, MSE allows users to search and retrieve relevant malware samples from multiple storage locations.

The key features of MSE include a powerful indexing system that organizes patterns based on attributes such as data hashes, behavioral patterns, and similar numbers. Machine learning algorithms automatically separate samples into different groups, making them easier to analyze quickly.

MSE application has a userfriendly interface that facilitates interaction between cybersecurity experts, researchers and analysts. Users can perform complex queries, see social patterns, and access detailed information.

Through observational analysis and realworld research data, we demonstrate, once and for all, the effectiveness of MSE in accelerating malware analysis and improving threat detection. Its scalability and adaptability make it suitable for integration into existing cybersecurity workflows, marking progress in the fight against malware threats.

**Keywords—**Malware, Search Engine, Research, Indexing, Malware Analysis, Cuckoo Sandbox, Machine Learning

## I. INTRODUCTION

Today's digital environment is increasingly affected by the spread of malware, which poses a major risk to network security. Detecting and combating these threats requires the development of effective tools to detect, identify and distribute malware samples. In response to this urgent need, this article introduces the design and implementation of the Malware Search Engine (MSE). MSE allows users to search and retrieve malware samples from databases using advanced data mining, machine learning, and data analysis methods. This introduction sets the stage for exploring the key capabilities and benefits of MSE in addressing the evolving challenges of malware detection and mitigation in cybersecurity.

### A. Search Engine

A search engine is a powerful tool that helps users navigate and retrieve information from the vast internet. It acts as a gateway to the vast amount of knowledge available online,

allowing users to find relevant content quickly and efficiently. In the digital age, search engines play an important role in facilitating information discovery and dissemination, whether you are looking for answers to questions, researching topics of interest, or exploring new ideas.

#### a. Phases of a Search Engine:

- a. *Crawling*: The search engine systematically traverses the web, visiting web pages and following links to discover new content. Specialized programs called web crawlers or spiders are deployed to explore the internet, indexing the content of web pages they encounter.
- b. *Indexing*: Once web pages are crawled, the search engine organizes and stores the information gathered in a structured format known as an index. The index serves as a catalog of the web, containing information about the content, keywords, and metadata of each indexed page.
- c. *Ranking*: When a user enters a query into a search engine, the system retrieves relevant results from its index. In the ranking phase, the search engine evaluates the retrieved documents based on various factors to determine their relevance to the query.
- d. *Retrieval*: When the user searches for malware, the search engine returns detailed results after ranking the data according to relevancy. Important details like the hash ID, download links, and other information relevant to the particular malware that was questioned are included in these results. It is possible for users to look up hash IDs directly, which makes it easier for them to obtain malware files. Researchers and security experts alike can conduct effective investigation and analysis thanks to this direct access to malware files.

## II. MOTIVATION

The primary motivation for developing the Malware Search Engine (MSE) stems from an understanding of the difficulties encountered by students and researchers, particularly those new to the field of malware analysis. Comprehensive malware analysis frequently requires access to a variety of malware samples for study and experimentation. However, sourcing such samples from the vast expanse of the internet can be difficult and time-consuming, particularly for those unfamiliar with the complexities of cybersecurity landscapes. MSE aims to reduce the burdens associated with manual search efforts across disparate online sources by consolidating a diverse range of malware specimens into a single, easily navigable interface. This initiative not only reduces the potential overwhelm for newcomers to the field, but also fosters an environment conducive to efficient and structured learning in malware analysis. MSE enables users to tailor their search criteria, allowing them to find relevant malware samples based on their research requirements. This personalized approach to malware analysis improves efficiency and precision. MSE aims to streamline this process by providing a centralized platform from which students and researchers can easily access a diverse repository of malware samples.

## III. MALWARE SEARCH ENGINE (MSE)

### A. *Intuitive Interface*

MSE features a user-friendly interface designed for seamless interaction, enabling cybersecurity professionals, researchers, and analysts to navigate the system effortlessly.

### B. *Streamlined Querying*

**Streamlined Querying:** Users can perform complex searches with ease, leveraging MSE's streamlined querying capabilities to retrieve specific malware samples based on various attributes such as file hashes, behaviour patterns.

### C. *Comprehensive Reporting*

MSE provides detailed reports that are easily accessible to users, facilitating further analysis and decision-making. These reports offer comprehensive insights into malware samples, aiding cybersecurity efforts effectively and efficiently.

### D. *Real-time results*

Leveraging advanced indexing and retrieval techniques, the search engine provides real-time results, ensuring that users receive the most up-to-date information on malware samples. This instantaneous feedback enhances the user experience, enabling rapid decision-making and response to emerging cybersecurity threats.

## IV. LITERATURE REVIEW

The article [1] offers a groundbreaking exploration of the spyware risk online. Using web crawling methods, the researchers conducted an extensive and ongoing analysis of the internet, examining both executable files and regular web pages to detect harmful content. The researchers sought to measure the presence of spyware on the internet by analyzing

executable programs and scripted drive-by download attacks. The study, conducted in May and October 2005, uncovered the prevalence of spyware in executable content, with a significant number of domains hosting spyware, including dangerous features like Trojan downloaders and dialers. Certain website categories, like games and wallpaper sites, were pinpointed as having a greater amount of executable spyware. Additionally, the research examined drive-by download attack occurrences over a span of five months, demonstrating a notable decrease in the number of pages facilitating such attacks from May to October. The study identified pirate sites and celebrity sites as the most dangerous web areas for drive-by attacks. The findings offered important information on the changing landscape of spyware online and highlighted the continued prevalence of spyware despite the decrease observed.

The article [2] introduces a method for web crawling that focuses on retrieving web pages related to specific topics. It highlights the significance of link classification methods in assessing the relevance of web pages, rather than categorizing the pages that have been downloaded. Through the integration of a Naïve Bayes classifier with URL scoring optimization, the study seeks to improve system efficiency in determining which links to crawl. The literature about focused crawlers shows how web crawling techniques have evolved to efficiently find information on specific topics. Several studies have looked into using machine learning algorithms like the Naïve Bayes classifier to evaluate the relevance of web pages to particular subjects. Progress in link classification methods has also played a key role in enhancing the precision and effectiveness of focused crawling systems. Researchers have studied how preprocessing methods like tokenization, stemming, and stop word filtering can improve the way web page content is represented for better information retrieval. Additionally, using seed URLs and topic words as starting points for the crawling process is a widely adopted practice in focused crawling approaches. In the area of focused crawling, challenges arise from the requirement for ongoing updates to classifier training sets throughout the crawling process to account for changing web content. Moreover, efficiently managing substantial data amounts during crawling while upholding precision and recall in link classification presents a significant hurdle for researchers. The article enhances the current literature with a specialized crawling method that uses machine learning and URL optimization to enhance the selection of relevant links for crawling. By tackling significant obstacles and integrating new techniques, the research progresses intelligent focused crawling for effective web information retrieval.

The article [3] explores the automated examination of malware through Cuckoo Sandbox. It underscores the importance of threat scoring in responding to malware incidents and points out possible constraints in using threat scores for classification purposes. Cuckoo Sandbox simplifies malware analysis through automated behavioral analysis, hash comparisons, and integrated tools for identifying malware and compromise indicators. Users can submit suspicious files for analysis in a sandbox environment without manual intervention. Detailed reports are generated detailing the actions of the suspected malware, such as system operations, file changes, communication with external systems, and activity screenshots. This automation enhances analysis

efficiency and offers crucial insights into malicious file characteristics. Threat scoring is crucial in malware incident response as it prioritizes security incidents by severity levels. For organizations with small incident response teams handling numerous security incidents, automated tools such as Cuckoo Sandbox are essential for promptly triaging and addressing threats. This method assists in pinpointing the most severe threats for immediate elimination and continuous monitoring of unresolved risks. By assigning severity scores to observed malware actions, incident analysts can efficiently prioritize their responses, resulting in faster remediation times and improved threat comprehension for post-incident analysis. Using threat scores for classification may have limitations. The current methodology of threat scoring in Cuckoo Sandbox could be arbitrary, potentially resulting in misleading classifications. The authors suggest a more robust methodology that concentrates on analyzing malware sample behavior. By examining API calls, contextual behavior, statistical analysis, and threat intelligence, a dynamic and dependable threat scoring system could be created. This method could improve the precision of threat prioritization and attribution of malware families, offering increased assurance for information security incident response teams.

The article [4] explores a latent semantic indexing classifier that merges link analysis with text content for the retrieval and indexing of domain-specific web documents. It aims to address challenges in focused crawling, like the requirement of initial training data, all while upholding a strong recall/precision ratio. The study evaluates the effectiveness of this method against other established web information retrieval strategies. The writers contend that using a statistical method for processing text and links is more effective than a knowledge-based logical representation system for topic-specific issues, such as creating vertical search engines. They propose leveraging content and link-related characteristics for both the classifier and the distiller functions of a targeted web crawler. Furthermore, the article explores the application of advanced algorithms like HCLA, PR, and others within this framework. Additionally, the article discusses the significance of integrating textual and linking data to optimize URL sequencing, citing methods that expand on PageRank or HITS. It highlights joint probabilistic models such as Probabilistic LSI/Analysis (PLSI) and Probabilistic HITS (PHITS) as instances of these strategies.

## V. METHODOLOGY

The MSE takes a similar approach to that of a standard search engine, however, instead of storing the samples themselves, we are storing the links to those malware samples scattered across the internet.

This approach brings several advantages including saving storage space and reducing the risk of inadvertently distributing malware.

The implementation of all phases with respect to MSE will be discussed in the following section.

### A. Crawling

Web crawlers, an essential component of the digital world, are automated software programs designed to navigate the internet in a structured manner. Referred to as web spiders or

bots, these tools are crucial for search engines in collecting, organizing, and categorizing online content, thus enabling easy access and searchability. Their process involves systematically exploring web pages through hyperlinks, allowing them to build a comprehensive index of the internet that improves the accuracy and efficiency of search outcomes. One notable type of crawler is the focused web crawler, which is specialized in retrieving content related to specific topics. These targeted crawlers enhance the quality and precision of search results by utilizing sophisticated algorithms to assess and select web pages and links based on their relevance to the given subject matter.

A focused web crawler's architecture is specifically designed to process information that meets its predefined criteria, making the data acquisition process more efficient. This is accomplished through detailed page and link scoring techniques that evaluate content relevance and establish the crawler's navigation path, ensuring that only relevant information is retrieved and indexed. The first step in implementing a focused crawler is selecting seed URLs, which are crucial as they determine the starting point of the crawling process. These URLs are usually chosen based on their expected relevance to the topic being explored, such as malware in cybersecurity applications. This selection is vital as it forms the basis for the crawler's ability to efficiently gather pertinent data.

Efficient management of seed URLs is vital for enhancing the crawling process. The seed URLs manager arranges and ranks these links to increase coverage and reduce duplicates, guaranteeing that the crawler consistently explores new and pertinent web areas. After a URL is handled, the data retrieved is parsed—an essential stage where crucial information is extracted, and fresh, relevant links are discovered. The parser carefully analyzes webpage content to sift and capture links that meet specific criteria relevant to the crawler's focus, such as potential malware indicators in cybersecurity applications.

Focused web crawlers are advanced tools that improve search engine efficiency and address specialized needs by collecting and filtering data relevant to specific queries. Their strategic process, from selecting initial URLs to analyzing content thoroughly, highlights their significance in the contemporary digital realm, where timely and relevant information retrieval is crucial.

### B. Indexing

During the indexing phase, we focus on collecting essential data from the web crawler in order to efficiently populate our database. Given the ethical considerations and limitations of obtaining detailed malware information via static analysis, we prioritize gathering fundamental attributes required for identification and retrieval.

At this point, our crawler extracts three important pieces of information: the malware's MD5 hash, its name, and the corresponding download link. Strategically indexing frequently queried fields like MD5 hashes and malware names in MongoDB significantly boosts query performance in the Malware Search Engine. By doing so, MongoDB can

swiftly locate and retrieve relevant malware samples based on user search criteria, enhancing overall efficiency and user experience. Additionally, indexing other attributes such as file types and timestamps further accelerates query processing, ensuring quick and precise retrieval of malware samples. While this dataset provides a foundation for indexing, it lacks critical details such as architecture and target operating systems, which are required for informed analysis and selection.

To close this gap and improve the comprehensiveness of our database, we use the Cuckoo Sandbox. By incorporating Cuckoo Sandbox into our workflow, we can automate malware analysis and gain valuable insights into the executables' characteristics.

We perform dynamic analysis of malware samples using Cuckoo Sandbox, which allows us to extract additional information such as architecture and operating system compatibility. This dynamic approach not only enriches our database, but also allows end users to make informed decisions based on their specific research or operational requirements. [5]

To supplement our indexing phase, in which we collect critical data such as MD5 hashes, malware names, and download links, we have implemented a Malware Classification Model to address the challenge of classifying malware. This model automates categorization by examining attributes like file hashes and behavioral patterns. By integrating this model, we improve the functionality of our Malware Search Engine, allowing for more precise cybersecurity analysis and research. Further information about the Malware Classification Model and the Cuckoo Sandbox will be provided in subsequent sections of our paper.

### C. Ranking

In the ranking phase, our focus lies on assessing the inherent threat posed by malware rather than evaluating the credibility of its source, without distinguishing between hard and soft links. Leveraging the Cuckoo community, we utilize user-defined scores assigned to API calls detected during malware execution within the Cuckoo sandbox. These scores are based on the frequency of API calls, which collectively contribute to the generation of a threat score for each malware sample. This score is calculated by summing the severity levels of matched Cuckoo signatures and normalizing it to a scale of 5.0. The variability in scores reflects the diversity of signatures matched against behavioral data, ensuring a nuanced assessment of malware threat levels. [3]

### D. Retrieval

Python is used in our system to act as a bridge between the NoSQL database MongoDB and the user interface. Python is used to handle user requests, which may contain hash IDs, when they start looking for malware on our website. The associated malware data, including hash IDs, download links, and other pertinent information, is then retrieved using Python interacting with MongoDB. Additionally, we used Python to create an intuitive database page for our website. Users can simply browse the whole malware database on this page, giving them a more comprehensive picture of the

information provided. Users can easily browse the database and learn about the wide variety of malware entries that MongoDB has stored.

All things considered, the amalgamation of Python, MongoDB, and our website's user interface streamlines the malware investigation procedure while also improving user experience. Our system makes sure that data is retrieved and accessible quickly, regardless of whether users are searching for specific information or looking through our malware database in its whole.

## VI. ARCHITECTURE OF MSE

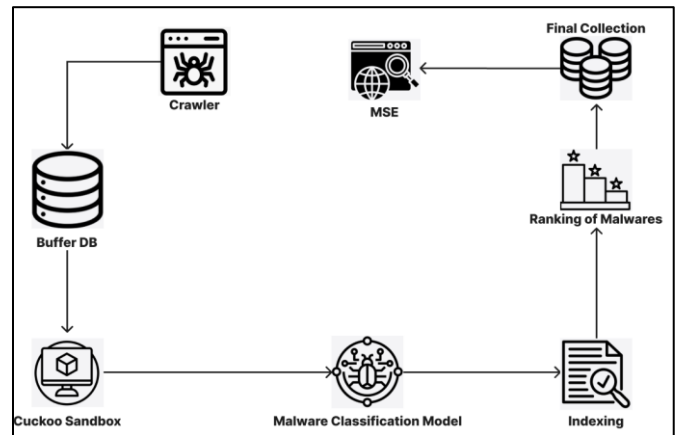


Figure 1. Depicts an overview of how MSE works.

Let's breakdown all the major tools we are using in the pipeline.

### A. MongoDB:

Using MongoDB offers us access to a cloud-hosted database, simplifying the data pipeline from crawler to collection. Its flexible schema design accommodates diverse metadata associated with malware samples, while its distributed architecture and scalability efficiently handle large data volumes. MongoDB's native support for replication and sharding ensures high availability and fault tolerance. Compared to traditional RDBMS, MongoDB's schema flexibility, scalability, performance, and agility make it an ideal choice for our Malware Search Engine, enabling efficient data management and robust storage and retrieval capabilities.

### B. Cuckoo-Sandbox:

The advantage of Cuckoo is its automation; no human interaction is required for behavioral analysis. A report is generated that describes the actions taken when the suspicious file was opened or executed, and all an analyst has to do is investigate the report. A detailed record of system operations, file manipulations, attempted communication with external systems, and even screenshots of the activity are created. All of these artifacts are required to understand the nature of the malicious file being analyzed.

### C. Python (3.8 /2.7)

Python (3.8 /2.7) plays a pivotal role in our project, serving as the primary programming language for automation. Its versatility and extensive libraries make it ideal for various tasks within the Malware Search Engine (MSE). Python facilitates seamless integration with MongoDB, enabling multiple operations on data within the collection. From automating the crawler function to downloading malware samples, preprocessing, post-processing, submitting samples for analysis to the Cuckoo Sandbox, and extracting relevant features from the generated JSON reports.

### D. Machine Learning

We employ the Extreme Gradient Boosting (XGBoost) classification algorithm, a machine learning model, for categorizing malware families within our system. This classification enables us to index each malware according to its respective family, enhancing organization and accessibility for users. By leveraging XGBoost, users can efficiently access specific types of malware families required for their analysis or research, further enhancing the functionality and usability of our Malware Search Engine.

## VII. IMPLEMENTATION

### A. Web Crawler

Web crawlers, also referred to as web spiders, web robots, or web bots, are computerized softwares applications that are designed to systematically and automatically search the world wide web in an structured way. By tracking hyperlinks from one web page to a different one, such crawlers explore online web pages, capturing data and classifying information for the search engine.

Web crawlers are crucial for retrieving data as they navigate the web in a organized way in order to collect, index and arrange the content. Search engines are able to locate, index and rank pages due to these automated computerized processes which provide quick access to relevant information. Web crawlers creates an index of web information by visiting the hyperlinks, thereby rendering accurate search results possible. They are essential tools for traversing the digital world as they support data extraction, information aggregation and quality assessment.

#### a) Focused Web Crawler Design and Architecture

A focused web crawler downloads only those pages that are relevant to a defined topic. Page scoring and link scoring methods are used by the focused crawler to figure out if a web page is about a specific topic. The appropriateness of a page is assessed via its score. The system of link scoring specifies which links should be crawled and in what sequence. For page scoring classification techniques are often employed. After downloading the pages, it utilizes the classification algorithm when deciding whether or not the content is about the defined topic. [6]

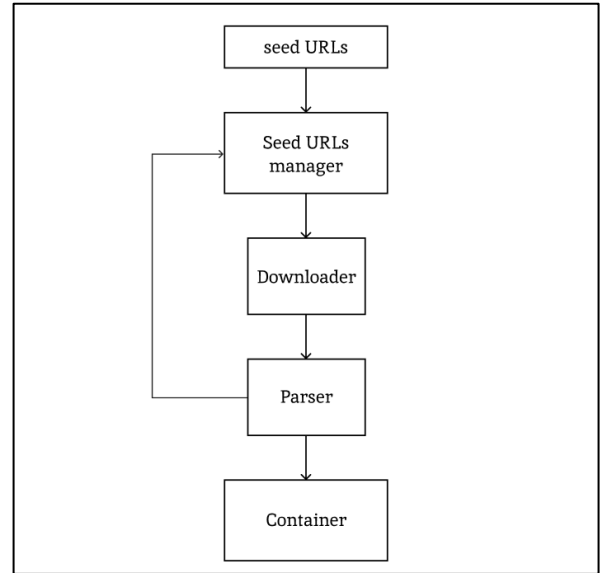


Figure 2. System architecture

#### b) Seed URL selection

The selection of seed URLs is vital because focused crawlers targeting malware samples have to enhance the degree of precision with which they retrieve web pages focused to a given topic. To determine the first n related pages, a malware related search keyword is first typed into an established search engine, for instance Google. These page URLs then subsequently make the initial batch of seed URLs that are added to the malware search engine's URLs manager. With the objective to enhance the effectiveness and appropriateness of the search process, this method ensures that the crawler begins with a highly appropriate list of starting points. Every web page is identified as a Node and an Edge indicates the link across any two pages. Additionally, all the seeds combined with each of their offspring pages forms a forest, with each seed page acting as the foundation of single forest tree, as seen in Fig. 2 [7]

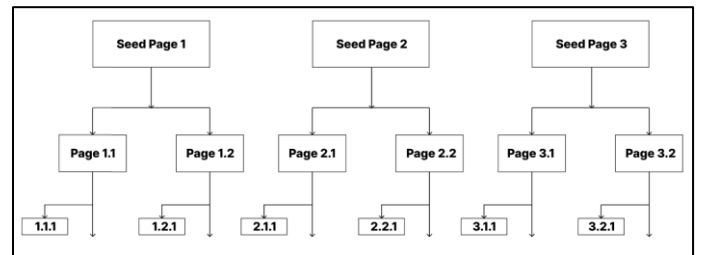


Figure 3. Web crawling forest

#### c) Seed URLs Manager

The seed URL is submitted to seed URLs manager at the beginning of the crawling process. URLs are effectively arranged and sorted for the crawling by the URLs manager, which maintains a priority queue. Repetitive visits to already crawled pages are prevented by the manager, which ensures that only unvisited URLs are enqueued and processed. This approach optimizes coverage, reduces duplicate attempts,

and directs the crawler towards new areas on the web. When the manager verifies that a new URL is present, the downloader receives the URL and begins the download process. The retrieved web page's data is fed into a customized parser, which eliminates the irregularities from the data and generates new URLs that we want in the seed URLs manager. The parsed data are stored in a local file system. The crawler then continues the process constantly until the seed URLs manager is emptied or a certain specified criterion are met.

d) Parser

The parser handles the crucial task of extracting appropriate data from webpage's HTML content. Its primary objective is to find potential hyperlinks to malware samples by thoroughly reviewing each hyperlinks and associated anchor text. Through systematic assessment, the parser identifies URLs or anchor text which includes terms that frequently correspond to malware in addition to file extensions that tend to be indicative of malicious programs, like ".exe" or ".dll". Links that fit these criteria are flagged for potential malware samples and saved for further examination and analysis. To put it in simple terms, the parser is an essential component of the web crawler's operation that makes it possible to identify and gather potentially harmful content from the internet.

B. Indexing

Indexing in our Malware Search Engine serves as the backbone of efficient data retrieval and organization. By strategically structuring and categorizing malware data, we enable users to swiftly locate and analyze relevant samples. Leveraging MongoDB, we store metadata such as MD5 hashes, malware names, and download links, facilitating seamless access to diverse malware specimens. Our indexing strategy prioritizes fields frequently queried by users, optimizing query performance and reducing response times. Additionally, techniques like sharding and replication ensure scalability and fault tolerance, guaranteeing uninterrupted access to the index even under heavy loads or hardware failures. Through meticulous schema design and optimization, our indexing process lays the groundwork for a robust and user-friendly platform for malware analysis and research.

a) Data Collection:

During the indexing phase, our focus is primarily on capturing essential metadata associated with malware samples encountered during the crawling process. Initially, we index key attributes such as the repository source, download link, and MD5 hash of the malware files. By prioritizing these fundamental identifiers, we lay the foundation for efficient retrieval and organization of malware data within our system. This streamlined approach allows us to quickly populate our index with valuable information essential for subsequent analysis and research activities.

```

{
  "_id": {
    "id": "65f21d607f2eb389e490cf1"
  },
  "repository": "ytsif/thezoo",
  "file_name": "PlasmaHTTP.zip",
  "download_link": "https://raw.githubusercontent.com/ytsif/thezoo/master/malware/Source/Original/PlasmaHTTP/PlasmaHTTP.zip",
  "md5_hash": "c2e955a41f48c40bfc9c537790bae PlasmaHTTP.zip"
}
```

Figure 4. JSON Schema of Indexed Data

b) Data Augmentation – Phase 1 :

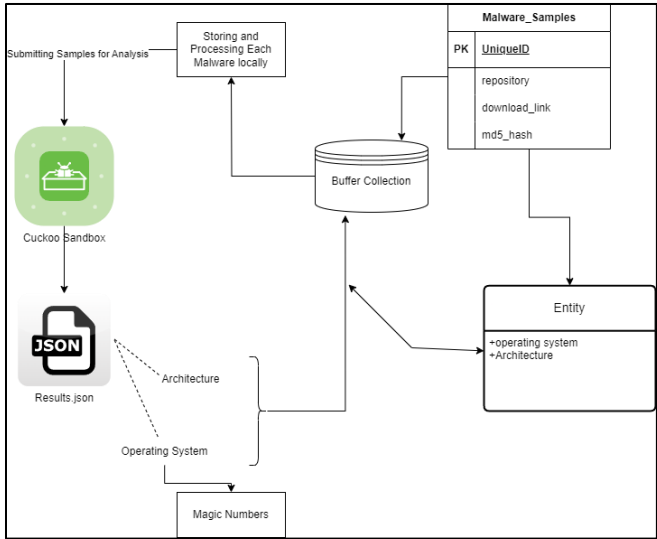


Figure 5. Phase 1 Architecture

This phase begins with a meticulous download of all collected malware samples into a controlled virtual environment, ensuring a secure and isolated space for subsequent processing. Each malware sample goes through a series of complex procedures designed to improve its quality and enrich its metadata.

One of the primary tasks during processing is unzipping the malware files, which frequently requires the use of specific passwords and extraction algorithms to access their contents. This step ensures that all relevant data from the files is extracted and available for further analysis. Furthermore, we meticulously sift through the extracted files to identify and remove any unwanted artifacts, such as .gitignore or MANIFEST files, as well as non-executable files, which could introduce noise into our dataset. This careful curation process helps to maintain the dataset's integrity and relevance, allowing for more accurate analysis results.



Figure 6. Stored Malwares inside Virtualenv

In addition to file extraction and cleanup, our automation scripts can handle a variety of processing challenges, such as



special characters and whitespace in file names. These scripts are critical in ensuring that all files are uniquely identified and properly processed, laying the groundwork for seamless integration into our pipeline.

To enhance the richness and utility of our dataset, we use the `os.walk` function and the `python-magic` library. This allows us to navigate through the directories containing malware samples, capturing file paths and corresponding magic numbers. By storing this information in a separate database file (`.db`), we add another layer of differentiation between executables and non-executables in our dataset. This meticulous categorization allows for more refined analysis and increases the accuracy of subsequent processing steps.

The integration of the Cuckoo Sandbox, which is well-known for its robust API endpoint support, is critical to our data augmentation strategy. Using this feature, a Python script manages the submission of each binary executable to the sandbox for dynamic analysis. This analysis includes a thorough examination of the malware's behavior and characteristics, revealing important details about its architecture and target environment. [5]

When we receive the `response.json` from the sandbox API, we extract useful information such as the architecture of each analyzed malware sample. This information is then seamlessly integrated into the malware's respective document, as identified by its MD5 hash field. This meticulous updating process ensures that our dataset is current and reflects the most recent analysis results, providing users with timely and actionable insights.

Central to our data augmentation strategy is the integration of the Cuckoo Sandbox, renowned for its robust API endpoint support. Leveraging this functionality, a Python script orchestrates the submission of each binary executable to the sandbox for dynamic analysis. This analysis encompasses a comprehensive examination of the malware's behavior and characteristics, yielding valuable insights into its architecture and target environment.

Upon receiving the `response.json` from the sandbox API, we extract pertinent information such as the architecture for each analyzed malware sample. This information is then seamlessly integrated into the respective document of the malware, identified by its MD5 hash field. This meticulous updating process ensures that our dataset remains current and reflective of the latest analysis results, empowering users with up-to-date and actionable insights. [8]

### C. Ranking

In the ranking phase of the malware search engine research paper, the emphasis shifts to evaluating the inherent threat posed by malware rather than assessing the credibility of its source. This means that the researchers are primarily interested in determining how dangerous a specific malware sample is, regardless of where it originated.

When malware executes, it communicates with the system and other software using Application Programming Interface (API) calls. These calls enable the malware to access files, change system settings, communicate across networks, and

execute commands. Each API call represents a specific functionality that the malware can use to achieve its goals, such as stealing data, spreading itself, or engaging in other malicious activities.

For example, let's look at the below example,

```
{
  "families": [],
  "description": "Queries for the computername",
  "severity": 1,
  "ttp": {},
  "markcount": 14,
  "references": [],
  "marks": [
    {
      "call": {
        "category": "misc",
        "status": 1,
        "stacktrace": [],
        "api": "GetComputerNameA",
        "return_value": 1,
        "arguments": {
          "computer_name": "NONE"
        },
        "time": 1715123111.187375,
        "tid": 2652,
        "flags": {}
      },
      "pid": 2648,
      "type": "call",
      "cid": 4481
    }
  ]
}
```

Figure 7. Sample of “GetComputerNameW” API call with a severity score of 1

Within the Cuckoo sandbox, cuckoo community has assigned severity scores to various API calls based on their potential threat level. These scores are assigned based on factors such as API call frequency and action type. For example, accessing sensitive system files may have a higher severity score than less critical operations. These community-defined severity scores are then combined to calculate a threat score for each malware sample. This threat score quantifies how dangerous the malware is based on the severity of its observed behaviour. To ensure consistency and comparability across samples, the threat score is normalized to 5.0. The threat score is calculated by summing the severity levels of all matched Cuckoo signatures associated with the malware's behaviour. Cuckoo signatures are patterns or indicators of malicious activity identified during the analysis. By incorporating these signatures into the calculation, the threat score becomes a more complete representation of the malware's potential impact. The variability in scores reflects the variety of signatures matched to behavioural data. This means that different malware samples may receive different threat scores, indicating the nuances of their malicious behaviour. Overall, this approach enables a more informed and precise assessment of malware threat levels, which aids in the detection and prioritization of security risks.

### D. Retrieval

#### Backend (Flask Application - `app.py`):

**Framework:** Our backend, which is based on the Flask web application framework, is made up of routes that can handle different HTTP requests, including GET and POST.

**Database Integration:** Our malware-based search engine's implementation smoothly combines frontend and backend elements, each of which is made to provide a user-friendly interface and effective data retrieval.

The program stores and retrieves data from a MongoDB database. It uses the supplied URI to establish a connection to a MongoDB Atlas cluster, and it uses the PyMongo module to communicate with the database.

Routes:

- **Renders the home page template (index.html).**  
**database:** Makes a connection to the MongoDB database, retrieves documents from the designated collection, and displays the retrieved documents by rendering the database.html template.
- **Search:** Responds to POST requests for document searches in the MongoDB collection using a query. Using several document fields, it conducts a compound text index search and provides the client with JSON data as the search results.
- **Error Handling:** incorporates an exception handler to manage exceptions from BadRequests. In the event that an invalid ObjectId is found, an error message is returned in a JSON response.
- **Custom JSON Encoder:** ObjectId serialization handling is implemented while returning JSON answers.

Frontend (HTML Templates and JavaScript):  
HTML Templates:

- **Database.html:** Stands in for the application's search screen. There is a form for users to submit search queries, and POST requests are sent to the server via JavaScript code that manages form submission.

MalwareX			
File Name	OS Type	Architecture & Platform	MD5 Hash
SAT Drop.zip	Windows	DOS batch file, ASCII text	93566f905c45d2513756344910e18
CryptLocker_10Sep2019.zip	Windows	PE32 executable (GUI) Intel x86_64, for MS Windows	2277895e3cd874ecf68322a7d9e5d
DC08 Yosemite.zip	Windows	DOS executable (COM)	5682a3c18f543c344c3d3d3d3d3d3d3d
EquationGroup_EquationGroup.zip	Windows	PE32 executable (GUI) Intel x86_64, for MS Windows	08a60cc70527a737074c33895932623
EquationGroup_Fanny.zip	Windows	PE32 executable (GUI) Intel x86_64, for MS Windows	4ec28b0c8b7c4387b0c303a9d4f9f562
EquationGroup_GROK.zip	Windows	PE32+ executable (native) x86_64, for MS Windows	286508a2b72c0d1867468ec3d1e5ef1
FancyBear_GermanParliament.zip	Windows	PE32+ executable (console) x86_64 (stripped to external PDB), for MS Windows	86c0dbf8eef769502c044147074d37
Friday the 13th:140.A.zip	Windows	COM executable for DOS	6d30c2b0d0e01c0d78f8b1c0a33b45
Friday the 13th:140.A.zip	Windows	COM executable for DOS	6d30c7035c40cb13c7f6e87f6c0c708
JS Lame.zip	Windows	Apple Desktop Services Store	688925e47f88f268895311e6fa3d15e
NOBEL.A.M.zip	Windows	PE32+ executable (DLL) (GUI) x86-64 (stripped to external PDB), for MS Windows	c0458a74a8a8976288917d35635a
NY10.S.zip	Windows	DOS/MBR boot sector, code offset 0x3c7, (EMSDOS MBR) rooted by Windows (MBR contains 2nd sector 2880) contains a MD5, serial#0717A, extension:16, serial number 3a11d1111, unbranded, FAT (12 bit), followed by FAT	34357a7c20c0d8b15444815319084d6d

- **Index.html:** depicts the application's home page and has a link that leads to the database search page (/databaseroute).



- **Styling:** For a consistent look and layout, both templates incorporate CSS styles. This includes general layout

tweaks as well as formatting for headings, forms, tables, and links. Table styles include background color alternating rows, cell padding, and border-collapse. In order to achieve a simple and intuitive design, form elements have defined widths, padding, borders, and hover effects. For a consistent look and layout, both templates incorporate CSS styles. This includes general layout tweaks as well as formatting for headings, forms, tables, and links. Table styles include background color alternating rows, cell padding, and border-collapse. In order to achieve a simple and intuitive design, form elements have defined widths, padding, borders, and hover effects.

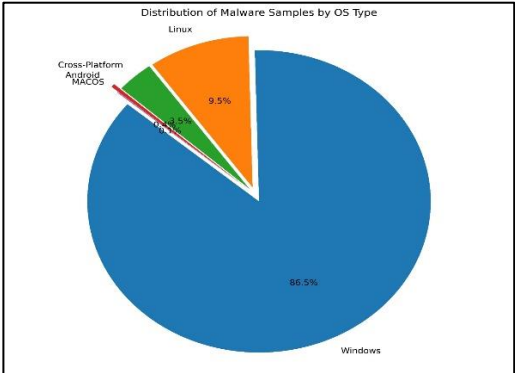
- **JavaScript (Search Functionality):** JavaScript code manages form submission and asynchronous search result retrieval in database.html. The JavaScript code dynamically creates an HTML table to display the search results after it receives the JSON input. A document from the search results is represented by each row in the table, which comprises columns for various fields (such as the file name, repository, and MD5 hash). For each "Download Link" column, clickable links are generated so that users can download the relevant files.



All things considered, this design combines strong backend functionality with an aesthetically pleasing frontend interface to offer consumers a smooth and captivating experience when searching and examining virus data.

VIII. RESULTS

A. Statistics

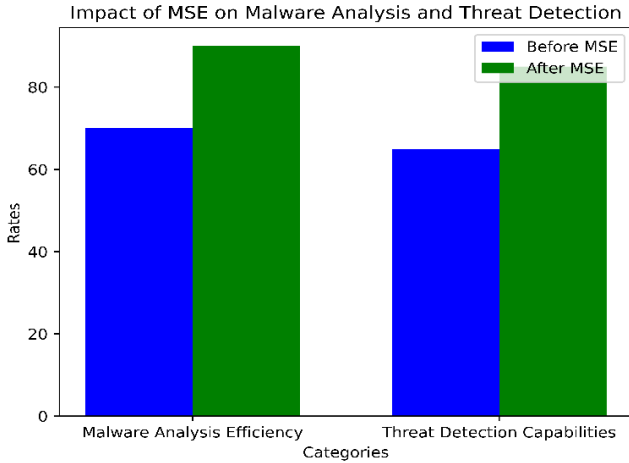




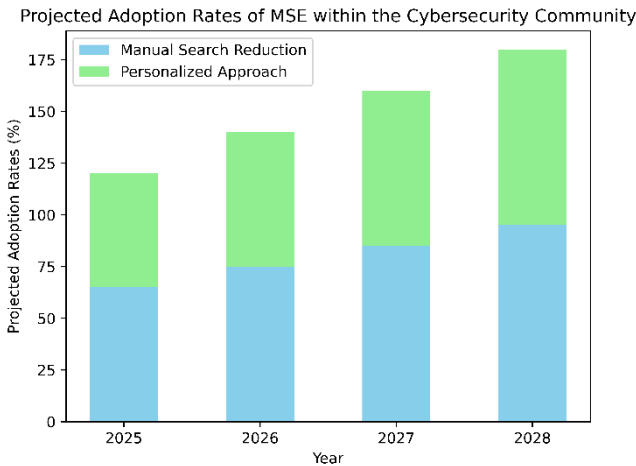
#### a) Number of Malware Samples Indexed or Stored:

The MSE consolidates a wide range of malware specimens into a single interface, making it easy for users to navigate. The system stores metadata like MD5 hashes, malware names, and download links to make retrieval and analysis easier. The indexing strategy prioritizes frequently queried fields, which improves query performance and reduces response time.

#### b) Impact on Malware Analysis Efficiency and Threat Detection:



The MSE streamlines the malware acquisition process, allowing security professionals and researchers to conduct efficient investigations and analyses. The MSE improves cybersecurity efforts by providing detailed reports and real-time results, as well as increasing malware analysis efficiency and threat detection capability. Adoption Rates within the Cybersecurity Community:



The MSE aims to reduce the burdens associated with manual search efforts across disparate online sources, while also creating an environment conducive to efficient and structured learning in malware analysis. The system's personalized approach to malware analysis improves efficiency and precision, appealing to users within the cybersecurity community.

## IX. CONCLUSION AND FUTURE RESEARCH

We presented the Malware Search Engine (MSE) in this paper as a novel tool intended to tackle the urgent problems in cybersecurity, specifically in the areas of malware analysis, detection, and classification. With the help of sophisticated machine learning, data analysis, and information retrieval techniques, MSE provides a robust indexing system that groups malware patterns according to a variety of criteria, including similar numbers, data hashes, and behavioral patterns. By using machine learning algorithms, MSE automatically divides samples into various groups, making analysis quicker and more effective.

There are encouraging directions for further improving MSE in the future. Expanding MSE's functionality would enable more thorough analysis across various operating systems by adding a specialized Malware Analysis Sandbox for malware based on Linux, macOS, and Android. Furthermore, by including a well-trained supervised machine learning model for malware classification during the data augmentation phase, end users may be able to access additional filters, meeting their research needs and improving the quality of the data.

In conclusion, the effectiveness of MSE in accelerating malware analysis and improving threat detection is demonstrated through observational analysis and real-world research data. Its comprehensive approach, coupled with its scalability and adaptability, underscores its significance as a valuable asset in the ongoing battle against malware threats.

## X. REFERENCES

- [1] T. B. S. D. G. a. H. M. L. Alexander Moshchuk, "A Crawler-based Study of Spyware on the Web," p. 17, 2006.
- [2] ". F. C. L. w. L. t. c. D. o. C. E. p. 5. 2. [1] D. Taylan
- [3] M. F. A. S. S. Aaron Walker, "Cuckoo's Malware Threat Scoring and Classification: Friend or Foe?," p. 7, 2019.
- [4] C. K. I. P. G. Almpandis, "Combining text and link analysis for focused crawling—An application for vertical search engines," p. 24, 2007.
- [5] S. Jamalpur, Y. S. Navya, P. Raja, G. Tagore and G. R. K. Rao, "Dynamic Malware Analysis Using Cuckoo Sandbox," *IEEE*, 2018.
- [6] D. Taylan, "Intelligent Focused Crawler: Learning which Links to crawl," *Department of Computer Engineering*, p. 5, 2011.
- [7] S. H. N. T. H. Z. F. L. L. W. Yongbin Yu, "A Survey about Algorithms Utilized by Focused Web Crawler," *JOURNAL OF ELECTRONIC SCIENCE AND TECHNOLOGY*, p. 11, 2017.
- [8] S. S. Darshan, M. A. Kumara and C. Jaidhar, "Windows malware detection based on cuckoo sandbox generated report using machine learning algorithm," *IEEE*, 2016.
- [9] J. Duckett, *HTML & CSS: Design and Build Websites*, 2022.
- [10] S. Krug, *Don't Make Me Think: A Common Sense Approach to Web Usability*, 2022.
- [11] J. VanderPlas, *Python Data Science Handbook*, 2022.
- [12] M. Haverbeke, *Eloquent JavaScript: A Modern Introduction to Programming*, 2022.
- [13] R. C. Martin, *Clean Code: A Handbook of Agile Software Craftsmanship*, 2022.
- [14] J. J. Garrett, *The Elements of User Experience: User-Centered Design for the Web and Beyond*, 2022.