# Signature-based Virus Detection for Windows Operating Systems Using Java and PostgreSQL

Eliyas Sala

## I. INTRODUCTION

Did you know that Windows is the most used operating system? In fact, 76.12% of computer consumers use Windows operating system in 2021 [4]. Since it has existed for such a long time compared to others, hackers and modern terrorists have found loopholes to manipulate computer systems running this operating system. They have created different malwares, like viruses, that are often undetectable by antivirus softwares [2]. The data given above also motivates individuals who seek business opportunities, like myself, to develop a software to secure this system. Although different antivirus softwares exist, most of them are ineffective and lack the consumers' trust. This research briefly discusses how these programs work and their tested performances. Then I propose my software-based solution with different development tools.

.

## II. RELATED WORK

Sukwong et al. examines the most used AV softwares and the algorithmic implementation behind them. He/she first provides the distinction between Signature-based detection vs. Behavior-based detection [5]. Both are ways to detect virus and every AV product doesn't necessarily use the same method to do that. In Signature-based detection, AV product scans a file and assigns a unique identification to that file. For example, it could use hashing algorithms like MD5 to assign value. Then it evaluates based on patterns of that hashing by comparing it to a remote database containing viral characteristics. If a match exists, then that file is indeed virus infected. Behavior-based detection, another common implementation, analyzes the behaviors of that specific file instead of accuracy matching. It is important to note that there are many more other implementations that different AV softwares use, but we will focus on these two.

Sukwong et al. continues to discuss and points out the ineffectiveness of AV softwares. Based on an experiment conducted to test the responsiveness of selected AV softwares, their activities were monitored when intentionally given infected files. The softwares in the experiment were the following: Avast, Kaspersky, McAfee, Norton, Symantec, and Trend Micro. Avast performed the best by detecting 62.15% of malware. However, they couldn't fully detect all of them. This clearly shows we shouldn't completely rely on our AV softwares to keep our windows computers secure.
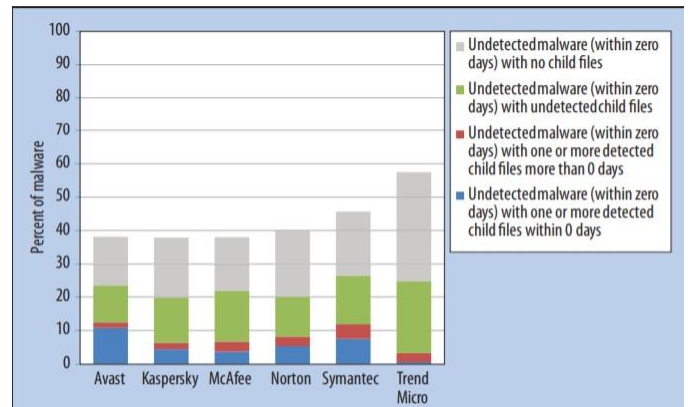


*Figure 1 shows comparisons of different AV softwares*

As you can see above, 42.83 percent of the malware remained undetected on first time. Among these undetected virus-infected files, 53.83 percent produced one or more malicious other files. Additionally, we found that 23.49 percent of these undetected malware files with malicious child files had one or more child files detected within zero days, but 64.36 percent weren't detected. This data tells us that there are limitations to these AV softwares.

## III. IMPLEMENTATION

To solve the problem of enhancing virus detection capability on Windows operating systems, it is important to analyze the different types of viruses that exist and learn their behaviors. Based on the statistics provided on the global market share of desktop operating systems, Windows has consistently dominated the market for years in staggering amounts. Other well-known operating systems, like macOS, are far from reaching this success. However, this much of consumer usage has created vulnerability for Windows operating systems [4]. Unlike most antivirus software, my software application runs detection-only and does not have a removal feature. It builds upon a Signature-based detection system.
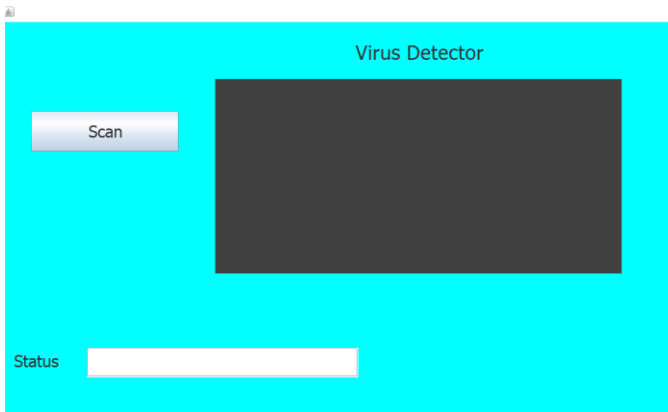
*Figure 2 shows my software*

I decided which programming langauge to use based on feasibility. I choose Java (my most proficient langauge) and I designed the user interface using Swing API. I created event listeners for the components as well. I imported Java's security class to use MessageDigest class to access it and generate MD5 hashing by scanning the files. The data structures I used include File object to store files as arrays and an arraylist to store MD5 strings. I also used a Scanner object to scan the files and store them as strings for manipulation.
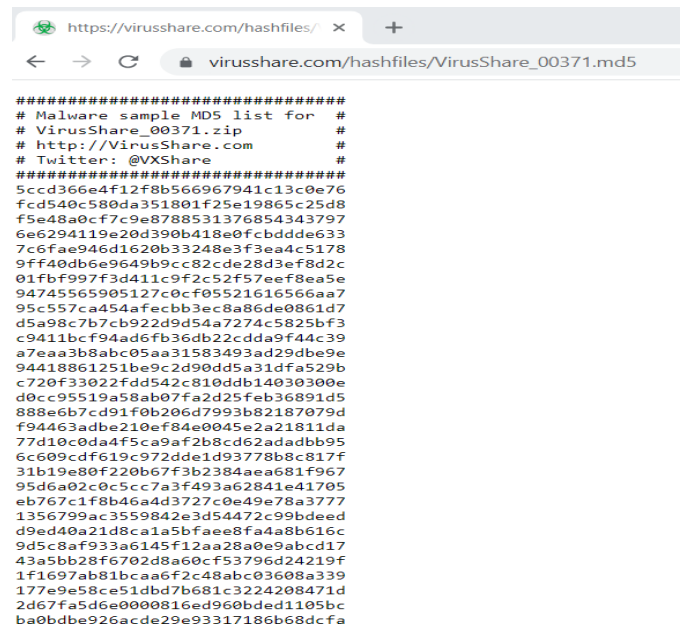
Given a directory path to scan for viruses, my software will first use a hashing algorithm to compute the MD5 hashing for every file in that directory. Then it matches strings collections of stored MD5 data in a PostgreSQL database containing large data of the existing viruses and their MD5 classification signatures [3]. After cross-referencing the MD5 generated from scanning the files against our data of known viruses' MD5, my application decides whether is it safe or virus infected. If infected, a pop-up message should appear before the users warning them.

IV.                          WORK

*Overall goal*

Deciding between Microsoft Visual Studio to develop a C/C# based program and Java-based desktop application. Since I have the most experience with Java, I chose to use Eclipse IDE to write my code and develop my program.



*Figure 3 shows java libraries and data structures I used*



*Figure 4 shows the data source I used for identifying viruses*

I used a real-world dataset from VirusShare.com. First, I stored the MD5 data in text files and then, I used Heroku to create a PostgreSQL database and added database driver to eclipse. Then I wrote a class for database connection and data insertion. Every time my software runs, it should either return "clean" or "virus detected". When users click "scan" button, it runs the scanning process to scan for viruses. The backbone of my application is the MD5 hashing function and Checksum functions that I wrote. I created a byte array to read data in chunks. Note that "checksum" function as a parameter of previous function's data. I incorporated multiple try and catch blocks since it is a good design pattern and conventional way to write a good code.

## V. CONCLUSION

This research-based project explores the different virus detection systems for Windows operating system. Then it discusses a signature-based scanning system with its implementation. Although my database only contains 10,000 rows of data due to resource limitations, the application is fully scalable and will continue to undergo refinement until my desired goals are met and optimal solutions are achieved.

## REFERENCES

[1] Al-Asli, Mohammed et al. "Review of Signature-based Techniques in Antivirus Products." (2019).

[2] Garba, Faisal A. et al. "Evaluating the State of the Art Antivirus Evasion Tools on Windows and Android Platform." (2019).

[3] Sahoo, Abhaya, et al. "Signature based Malware Detection for Unstructured Data in Hadoop." (2014).

[4] "Desktop operating system market share worldwide," StatCounter Global Stats. [Online]. Available: https://gs.statcounter.com/os-market-share/desktop/worldwide. [Accessed: 05-Sep-2021].

[5] O. Sukwong and H. S. Kim, "Commercial Antivirus Software Effectiveness: An Empirical Study," IEEE Computer Society, pp. 63-70, 2011.