

# How to Write YARA Rules That Minimize False Positives



Written by Ari Eitan - 12 May 2022

## Generate Advanced YARA Rules Based on Code Reuse

Incorporating YARA into daily security operations can accelerate incident response time, classify malware, empower threat intelligence and improve detection capabilities by creating custom signatures.

While YARA is a popular tool for SOC and IR teams, the main challenge is deciding what to base your YARA rules on for maximum effectiveness. In this post, we will explain how identifying code reuse between malicious files can be used to automatically develop advanced YARA rules to increase the accuracy of malware detection, reduce false positives, and improve threat hunting capabilities.

## What Are YARA Rules?

YARA, short for “yet another recursive acronym,” is a tool used in malware detection and classification. Malware researchers leverage YARA to create descriptions of malware families based on textual or binary patterns. Each description, or signature, is a YARA “rule” consisting of a set of strings and a boolean expression to determine its logic. When written effectively, YARA rules identify commonalities in malware and classify malicious files to other forms of malware that display similar patterns. More advanced YARA rules can be used to find additional variants of malware and hunt for samples.

([Click here to jump down to examples of code-based YARA rules.](#))

While most security vendors use their own signature mechanisms, YARA is an open standard tool which can be used with many platforms and applied to different use cases. For example, YARA enables an organization to create detections of its own, which is relevant for identifying targeted malware that generic or classic signatures often struggle to detect. This poses a significant advantage for security teams since they do not need to wait for signature updates from their security vendors. Security teams can create personalized, sophisticated YARA signatures even for targeted attacks and custom created malware.

## The Challenges in Writing Effective YARA Rules

The vast majority of YARA rules available today are simple-based rules focused on string reuse found in a malware's binary. Strings can include a log message or hard-coded user agent, which are criteria not guaranteed to be unique. Therefore they can result in false positives and can be easily replaced or encrypted by the adversary to avoid detection.

The most effective YARA rules are designed to achieve high detection and classification rates while reducing the number of false positives. Researchers must select the right textual or binary patterns to optimize detection results and accurately classify a file to its respective malware family. This is difficult because files share hundreds and even thousands of strings and it is imperative that rules are not too generic nor too specific. For example, if a researcher is defining a YARA rule based on a file that contains an embedded library, and the rule is based on the generic library alone — which is not a malicious piece of code in itself — the researcher will receive a hit on every single file that uses this library, in other words increasing the number of false positives generated.

The ideal signature will have the right balance: broad enough to identify many variants of the malware but specific enough to avoid false positives.

Here's an example:

```
101100010011100101010101  
0110101011001011000011  
1000011011110001101011  
1001100110101011101010  
0110110101100101110110  
0110110101100101110110  
1110101001110011011010  
1000011011110001101011  
1001100110101011101010  
0110110101100101110110  
01110110011100010110  
1110101001110011011010  
1000011011110001101011  
0110101011001011000011  
1110101001110011011010
```

The above example represents a file that has been disassembled into tiny pieces of binary code, or genes. If a YARA rule is created based on the file's trusted code (green), many false positives will be generated.

On the other hand, ensuring a YARA signature is not too specific is also important. Taking a piece of code as it is, for example, and developing a rule that contains all of the code's addresses will make the YARA rule very specific to this particular file only. The rule will not likely generate hits for other or future variants of the malware that contain different values.

### Writing Advanced YARA Rules is Difficult to Scale

More complex YARA rules can be created by incorporating features such as wild cards (a type of hexadecimal string), case-insensitive strings, and regular expressions. While advanced YARA rules can be powerful instruments for detecting malware, writing them requires a high degree of technical ability and an understanding of YARA, and can be a time-consuming process. Writing advanced YARA rules and [malware analysis](#) in general is difficult to scale at high volumes, especially since not every organization has access to a team of highly-skilled researchers or reverse engineers. How can a security team create advanced YARA rules and still achieve automation, especially if the organization is faced with a high volume of alerts? The answer lies in studying patterns in code reuse.

### Defining Code Reuse

Almost every software or malware today contains previously written code. Developers of trusted applications will employ code reuse to make their work more efficient and to bring tools to market faster. The same approach applies for malware authors. As attackers write more malware they will establish code patterns. For defenders, this provides an opportunity for attribution and identifying threat actor capabilities.

Genetic analysis of threats breaks down a given file into thousands of tiny fragments of code, or genes. Identifying code that was used in previous attacks can provide critical insights for security teams, including:

- Accurately determining the intent of the software. For example, is the file trusted or malicious?
- Classifying a malicious file to its relevant malware family.
- Uncovering the level of sophistication and potential risk of the threat. For example, are you dealing with a common banking trojan, a sophisticated APT or a nation-state sponsored attack? The answer will shape your response.
- Making attribution to the threat actor responsible for creating the malware. Not only can studying patterns in code reuse identify the origin of any given file, it can highlight unique, never-before-seen code, which can detect new threats that have been written from scratch.

### Generate Automatic, Advanced YARA Rules with Code Reuse

Just as attackers reuse code to deploy new malware, defenders can identify patterns in code reuse to create advanced YARA rules. Here's how:

- By identifying a malicious file's unique binary code, a signature can be produced for detecting only those genes. This will enable detection of the exact same malware with high confidence and a low false positive rate.
- Detecting samples of the same variants (the same opcodes with different addresses) can be replaced with wildcards.

- Detecting different variants of the same malware family or campaign (same baseline of code but with different functionalities or capabilities) to look for a partial match.
  - Detecting future variants that the attacker may release based on code reuse.

Let's look at that example again:

```
1011000100111001010101  
011010101011001011000011  
1000011011110001101011  
1001100110101011101010  
0110110101100101110110  
0110110101100101110110  
1110101001110011011010  
1000011011110001101011  
1001100110101011101010  
0110110101100101110110  
0111011001111000010110  
1110101001110011011010  
1000011011110001101011  
01101010110010111000011  
1110101001111000110110
```

As demonstrated in the previous example, creating a YARA rule based on a file's trusted code will result in false positives. At the same time, developing a rule that contains all of the code's addresses will make the rule very specific to this particular file only. The rule will not likely generate hits for other or future variants of the malware that contain different values.

By identifying the malicious (red) and unique code (purple), a researcher can create a YARA rule that will achieve very accurate results. With a YARA rule based on code reuse, you can expect to:

- Receive hits on different hashes that contain the exact same code.
  - Receive hits on different hashes that contain some but not all of the original code. This is useful for identifying new variants that are similar to the malware that the signature is based on.

## Example YARA Rule for Emotet Malware

The following example demonstrates a code-based YARA rule produced for [Emotet](#), a common financial trojan. With the code-based YARA signature in place, any future sample or variant that reuses some aspects of Emotet's code will be detected.

Example YARA rule, generated by [Intezer](#) to detect Emotet.

The screenshot shows the Malicious platform interface. At the top, there's a navigation bar with tabs: Genetic Analysis (selected), TTPs, IOCs, Behavior, Detect & Hunt (Beta), and Extended Dynamic Execution. The main content area has several sections:

- Original File**: SHA256: 5f2aa0bb76749cce2c13c1049521a3e70389c773632245e1f915ca6c522d1402, VIRUSTOTAL Report (38 / 60 Detections). Status: non\_executable, excel.
- Dynamic Execution**: Powered by Cape. Shows tasks like explorer.exe | 1660, EXCEL.EXE | 2508, regsvr32.exe | 1528, and regsvr32.exe | 2940.
- Code clusters**: A table showing Related Families (114 genes) and Code Cluster (0x418 (89 Blocks)).
- Code (114)**: A detailed table of assembly code blocks.
- Strings (18)**: A list of strings found in the file.
- Capabilities**: A list of capabilities observed.

A red box highlights the "Generate YARA" button in the "Code (114)" section. The bottom right corner features a "Help" button with a question mark icon.

[5f2aa0bb76749cce2c13c1049521a3e70389c773632245e1f915ca6c522d1402](#))

## Example YARA Rule for WannaCry Variant2

This is an example of a YARA signature that was produced via code from a [WannaCry](#) sample. Deployed in May 2017, WannaCry was one of the largest, high profile ransomware attacks in history, infecting over 200,000 computers across 150 countries.

Example YARA rule for WannaCry.

Genetic analysis of WannaCry in Intezer Analyze (SHA256: [bf293bda73c5b4c1ec66561ad20d7e2bc6692d051282d35ce8b7b7020c753467](#))

## Generating YARA Rules with Intezer Analyze

Leveraging proprietary genetic code analysis technology, Intezer Analyze can generate automatic YARA signatures based on a file's code, enabling users to improve their threat hunting capabilities by detecting future variants of the malware. As highlighted above, code-based YARA signatures can effectively reduce false positives and save precious resources in the form of time and analyst efforts. This is particularly valuable for organizations dealing with a large volume of alerts.

You can also use [Detect & Hunt](#) to find file-based detection opportunities for creating new YARA rules. You can filter these detection opportunities by the artifact type, family, and verdict.

NanoCoreRAT sample, SHA256: [0689544cd227b93df10df4d99ea04270d0f9f76259aa34cb8f6707a02e70081d](#)

This feature also allows for some more advanced usages. It can be used in several scenarios, including:

1. The user can adjust the thresholds of the rule. The default value is 70% of the entire code but it is possible to modify the rule to be more specific, or more flexible.
2. The user can combine or split different YARA signatures to build stronger rules to better fit his or her needs.
3. The user can add to the automated rules with a string reuse feature to make the signature even more powerful for threat hunting.

## Key Takeaways for Writing Effective Rules

- Incorporating YARA into daily security operations can accelerate incident response time, classify malware, empower threat intelligence and improve detection capabilities by creating custom signatures.

- Traditional methods for writing effective YARA signatures have their challenges, including being too specific or generic with respect to a file's textual or binary patterns.
- In today's YARA landscape many signatures are based on string reuse. The challenge is to identify the "unique" strings, however even if this is achieved, these signatures are much less effective because strings can be easily manipulated, replaced or encrypted by the adversary to avoid detection.
- Identifying malicious code seen in previous threats can be used to generate more accurate YARA signatures for detecting future variants or new malware.
- The Intezer Analyze platform enables users to automatically generate YARA rules based on binary code, rather than simple strings. Code-based YARA rules are the most effective since they are tolerant to modifications and are more equipped to detect variants of the same threat.

[Create a free Intezer account here](#) — start analyzing files to generate custom YARA rules or download YARA rules for public malware samples in Intezer now.



Ari Eitan

Ari manages the team responsible for the genetic algorithm behind Intezer's code genome database. In his role as VP of Research, Eitan leads the company's malware hunting and investigation operations, analyzing threats and publishing information about new APTs. Eitan began his career as a security researcher for the Israeli Defense Force (IDF). He quickly became Head of the IDF's cyber incident response team (IDF CERT), honing his expertise in incident response, malware analysis, and reverse engineering. Eitan has presented his research at several government and information security events, including AVAR, BSidesTLV, CyberTech, Hack.lu, Hacktivity, Infosec, IP EXPO, Kaspersky SAS, and the Forum of Incident Response and Security Teams (FIRST).

[code reuse](#) [detection rules](#) [Emotet](#) [false positives](#) [Threat Detection](#) [WannaCry](#) [yara rules](#)