# How to Interpret MITRE ATT&CK® Evaluations

## Expert tips on how to assess security product capabilities using The MITRE Foundation's landmark research

## Executive Summary

When looking to improve their security posture, buyers will often look to third-party analysts to assess technologies and vendors. This is a good resource, but often their reports are built on survey data and as a result, vary from source to source. Furthermore, they often can't assess the actual effectiveness or performance of products in the face of living threats.

Thankfully, The MITRE Foundation, a not-for-profit entity, provides a publicly available cyber-defense test of endpoint security products called the ATT&CK Enterprise Evaluations. Because of the report's depth and lack of a stack rank and scoring system, people tend to gravitate toward easily consumed vendor-generated graphics and charts composed from cherry-picked data, which puts buyers at the mercy of misinformation. Readers must take the time to ingest and interpret the information to understand how each offering did and why, and whether each vendor's technology approach may or may not match your organization's security strategy needs.

Starting in 2019, the MITRE ATT&CK Evaluations have provided four rounds (sometimes called "phases") of detailed tests of the capabilities of endpoint security solutions by emulating real-world cyber campaigns and their techniques and tactics.

This paper will walk you through step by step to understand the fundamentals of the MITRE ATT&CK Evaluation so you can more easily navigate the results when making decisions on which vendors to evaluate to secure your endpoints and integrate into your security strategy.

## Background On the MITRE ATT&CK Evaluations

Starting in 2019, the MITRE ATT&CK Evaluations have provided four rounds (sometimes called "phases") of detailed tests of the capabilities of endpoint security solutions by emulating real-world cyber campaigns and their techniques and tactics. The MITRE Engenuity ATT&CK Evaluations are powerful because they are based on the MITRE ATT&CK framework, which is a robust knowledge base of adversarial techniques. It provides a breakdown and classification of offensive actions taken by attackers that can be used against particular platforms, such as Windows. Unlike prior work in this area, the focus isn't on the tools and malware that adversaries use but on how they interact with systems during an operation.

To provide context, the ATT&CK framework organizes techniques into a set of tactics (**what** the cybercriminal is attempting to do), each with specific techniques (**how** they try to do it). Each technique includes information that's relevant to defenders to help them understand the context surrounding events or artifacts generated by a technique in use. The relationship between tactics and techniques can be visualized in the ATT&CK Matrix, which spans 14 discrete techniques. The Matrix offers a robust and granular mapping of the activity of potentially utilized cyberattacks. Each area has seven or more tactics and spans from reconnaissance through impact.

The 2022 round of MITRE Engenuity ATT&CK tests focused on two threat actors, Wizard Spider and Sandworm. Wizard Spider is a financially motivated criminal group that has been conducting ransomware campaigns since August 2018 against a variety of organizations, ranging from major corporations to hospitals. Sandworm is a destructive Russian threat group that is known for carrying out notable attacks such as the 2015 and 2016 targeting of Ukrainian electrical companies and 2017's NotPetya attacks. These two threat strains were chosen based on their complexity, relevancy to the market, and how well MITRE Engenuity's staff can fittingly emulate the adversary.

In the latest evaluation, MITRE first ran the detection test to see what sub-techniques the endpoint security solution will detect (and present with or without context), followed by the protection test to see if or when it will block the attack. Through careful examination of the screenshots, you can gain a better understanding of the usability of the products and the manner of the protection it provides.

## The Detection Test

The 2022 test comprised 19 steps with multiple stages called either *sub-techniques* or *sub-steps*. The evaluations used six terms to express how the product performed for each test and noted the data source for the detection. Depending on a vendor's participation or detection abilities, you will see one of these six detection terms in order of value:
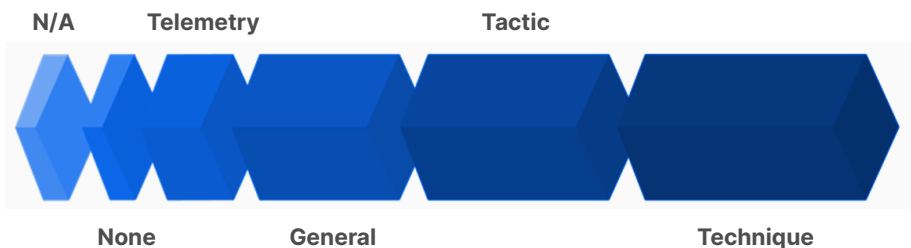


Figure 1: MITRE Detection Categories

- **Technique:** In short, the ideal outcome. The behavior was processed and designated malicious, as well as enriched with the technique or other notation about how it was performed by the attack.
- **Tactic:** The behavior was processed and designated malicious, as well as enriched with the tactic or other notation about why it was performed by the cyberattack.
- **General:** The behavior was processed and flagged, but without detail as to why (tactic) or how (technique) the action was performed.
- **Telemetry:** The behavior could be seen but was minimally processed.
- **None:** In short, the least ideal outcome. No data indicating the test behavior was detected could be seen within the product.
- **N/A:** Seen on the Linux test (sub-techniques 11.A.1 through 14.A.5) for those that did not participate in this portion (nearly a third of vendors). This is a neutral result; talk to these vendors about OS coverage if Linux protection is in scope.

More importantly, MITRE publishes robust information on each product evaluated, including:

1. Total number of detections tested and made (this is titled "Visibility") and would appear as x of y sub-techniques.

2. Total number of detections made with the MITRE technique noted (this is titled "Analytic Coverage") also provides a linked screenshot so that those interested can see the user experience firsthand.

3. A list of all sub-techniques and detections (including reference screenshots to show the user experience for each detection).

To start, go to the overview page, select a vendor, select "Wizard Spider + Sandworm (2022)," then scroll down to start with "Scenario 1" (Wizard Spider) or click on "Scenario 2" (Sandworm, which includes an optional Linux test). Upon scrolling down on a scenario, you will see the sub-techniques (e.g., 5.A.8) and their detection type (e.g., Technique).
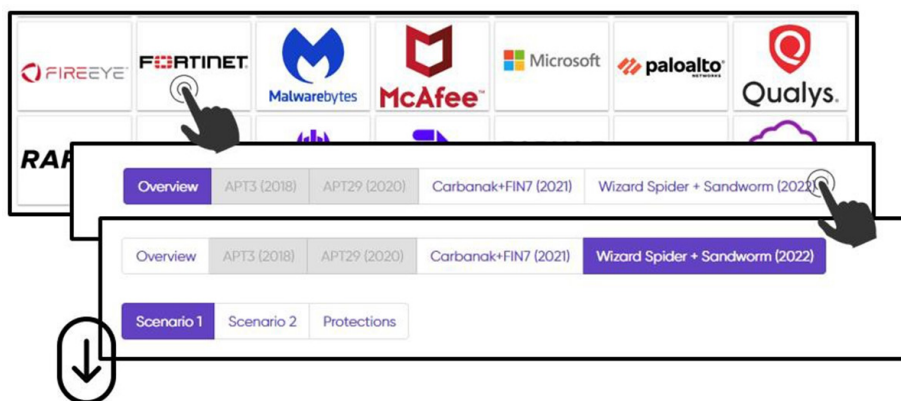


Figure 2: How to Navigate to the Detection Results

| 5.A.9 | **Tactic**<br>Discovery (TA0007)<br><br>**Technique**<br>Permission Groups Discovery (T1069) | **Tactic**<br><br>**Technique** (Configuration Change (Detection Logic)) |
|---|---|---|
| 6.A.1 | **Tactic**<br>Command and Control (TA0011)<br><br>**Technique**<br>Ingress Tool Transfer (T1105) | **Telemetry** |
| 6.A.2 | **Tactic**<br>Credential Access (TA0006)<br><br>**Technique**<br>Steal or Forge Kerberos Tickets (T1558)<br><br>**Subtechnique**<br>Steal or Forge Kerberos Tickets: Kerberoasting (T1558.003) | **General** (Delayed)<br><br>**Technique** |

Figure 3: A Sample of Detection Results From Scenario One

## Configuration Changes

As you go through the detection results in both scenarios, you will notice configuration changes. Some vendors will try to chart these vendor by vendor while claiming that these are delays in detection; within the milliseconds, an attack will initiate and potentially cause damage. At times they may be correct, but certainly not in all cases. MITRE will note when and why a configuration change occurs within their test. You could see a change in logic, which might be due to the client asking the strain to attack again to retest the results—others, a change in data source. These types of configuration changes introduce no latency. Also, MITRE will note when there is a delay. These are often due to waiting for an analyst and sandboxing results. These can introduce latency and is worthy of note due to allowing attacks to continue along the kill chain unchecked.

| 1.A.10 | **Tactic**<br>Command and Control (TA0011)<br><br>**Technique**<br>Application Layer Protocol (T1071)<br><br>**Subtechnique**<br>Application Layer Protocol: Web Protocols (T1071.001) | **None** (Configuration Change (Detection Logic))<br><br>**Telemetry** |
|---|---|---|
| 1.A.11 | **Tactic**<br>Command and Control (TA0011)<br><br>**Technique**<br>Encrypted Channel (T1573)<br><br>**Subtechnique**<br>Encrypted Channel: Symmetric Cryptography (T1573.001) | **Technique** (Configuration Change (Data Sources)) |
| 2.A.1 | **Tactic**<br>Persistence (TA0003)<br><br>**Technique**<br>Boot or Logon Autostart Execution (T1547)<br><br>**Subtechnique**<br>Boot or Logon Autostart Execution: Registry Run Keys / Startup Folder (T1547.001) | **Technique**<br><br>**Technique** (Delayed) |

Figure 4: Configuration Changes in Detection Results

## Analytic vs. Telemetry Coverage

After clicking on a vendor on the overview page, instead of proceeding to a vendor's 2022 test, take a look at the overview chart (see the middle screenshot in Figure 2 above or Figure 5 below). From here, you can see how one did in "Analytic Coverage" as opposed to "Telemetry Coverage." Analytic coverage is ideal as it is the easiest to understand because it calls out the specific activity detected and uses the emerging industry-standard attack lexicon to describe it. When compared to telemetry at the other end of the spectrum, the activity is only logged and can be found, but with effort and must be interpreted based on the vendor's syntax. As customers mature along with the industry, detections based on analytics are preferred for more accurate detections, especially when event data is correlated between multiple solutions. Note: When reviewing the overview section, telemetry detections can be tallied twice but duplications are scrubbed out in the "Visibility" box for the test in the latest round 4. (In previous rounds, multiple detections were counted cumulatively.)

Evaluation Summary
These are the evaluations that WithSecure has participated in:

| Evaluations | Analytic Coverage ⓘ | Telemetry Coverage ⓘ | Visibility ⓘ | Detection Count ⓘ |
|---|---|---|---|---|
| APT3 (2018) | 72 of 136 substeps | 107 of 136 ⓘ substeps | 122 of 136 substeps | 217 across 136 substeps |
| APT29 (2020) ⚪ Include MSSP | 90 of 134 substeps | 110 of 134 ⓘ substeps | 118 of 134 substeps | 224 across 134 substeps |
| Carbanak+FIN7 (2021) | 80 of 174 substeps | 137 of 174 ⓘ substeps | 152 of 174 substeps | 253 across 174 substeps |
| Wizard Spider + Sandworm (2022) | 66 of 109 substeps | 17 of 109 substeps | 83 of 109 substeps | ⊖ |

Figure 5: The Overview Results for a Vendor's MITRE Evaluation Results

## The Protection Test

After the completion of the detection test, the protection test commences to see if and where within the kill chain the attack is stopped. To navigate to this test, from the overview page, select a vendor, select "Wizard Spider + Sandworm (2022)," and then select "Protections."
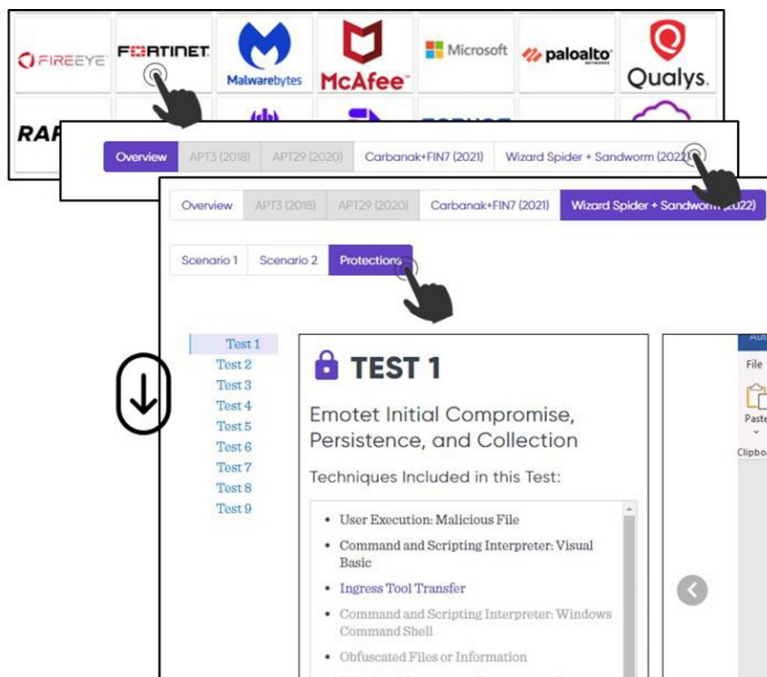


Figure 6: How to Navigate to the Protection Results

This set is much easier to diagnose the results, which only uses these three terms represented by different icons:

**Blocked (blue closed lock):** The technique was blocked, and the user was informed that the attack was unsuccessful.

**None (black open lock):** There was no evidence that the technique was blocked or otherwise unsuccessful because of the product.

**N/A (gray lock with a line drawn through):** This is reserved for those who did not participate in the Linux test. Those that did not participate in any protection tests will have a grayed-out "Protections" cell next to "Scenario 2."

In Figure 6 above, or in any protection test you are reviewing, you will notice a list of "Techniques Included in This Test." This will show you the sub-techniques and, if it was blocked, the step where the block occurred in blue text. Black text represents steps where the attack proceeded down the chain, and gray text is steps that would have occurred if the attack continued past the block or in the case of test seven (Linux) if the vendor did not participate. Although there are fewer terms in the protection test compared with the detection test (three terms vs. six), the interpretation is more complex because the timing and the method of blocking are important and can vary based on vendor and the test case.

At this juncture, you can compare various vendors' results to see where they blocked within the kill chain. Those that block based on signature, machine learning (ML) or other static analysis will block very early in the process, as is the nature of knowing the attack. If unknown, the attack may continue and cause damage in the end, so this is why many modern solutions continue to conduct behavior-based analysis and ideally protection. While it is good if attacks are blocked early, doing so can increase false positives. In contrast, blocking too late may expose the organization to a degree of risk even if the end objective is not achieved.
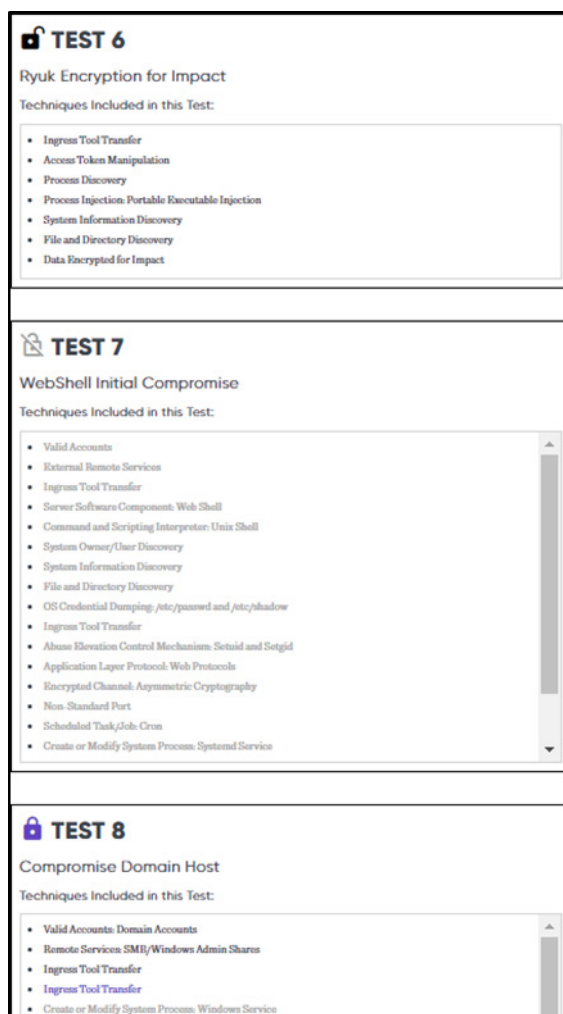


**TEST 6**

Ryuk Encryption for Impact

Techniques Included in this Test:

- Ingress Tool Transfer
- Access Token Manipulation
- Process Discovery
- Process Injection: Portable Executable Injection
- System Information Discovery
- File and Directory Discovery
- Data Encrypted for Impact

**TEST 7**

WebShell Initial Compromise

Techniques Included in this Test:

- Valid Accounts
- External Remote Services
- Ingress Tool Transfer
- Server Software Component: Web Shell
- Command and Scripting Interpreter: Unix Shell
- System Owner/User Discovery
- System Information Discovery
- File and Directory Discovery
- OS Credential Dumping: /etc/passwd and /etc/shadow
- Ingress Tool Transfer
- Abuse Elevation Control Mechanism: Setuid and Setgid
- Application Layer Protocol: Web Protocols
- Encrypted Channel: Asymmetric Cryptography
- Non-Standard Port
- Scheduled Task/Job: Cron
- Create or Modify System Process: Systemd Service

**TEST 8**

Compromise Domain Host

Techniques Included in this Test:

- Valid Accounts: Domain Accounts
- Remote Services: SMB/Windows Admin Shares
- Ingress Tool Transfer
- Ingress Tool Transfer
- Create or Modify System Process: Windows Service

Figure 7: Three Types of Protection Results

Suppose that test 1 (steps 1.A.1 through 3.A.5 in this test) was blocked at the first step 1.a.1. The result sounds great; the cyberattack was stopped at the earliest stage possible. But what if it was actually the user execution that was blocked? In that case, you would want to know the basis on which the user was prevented from accessing a file. Was there a high-confidence malicious indicator, or were policies set too aggressively? Alternatively, suppose blocking occurred at the end of step 3.a.5, which was email collection. In that case, the product stopped the intended data breach, but it did allow step 1.A.8, remote file copy, to occur, which means the attack had a malicious impact. Often, the impact of an attack can have several negative consequences.

In this particular case, arguably the safest time to block the attack to minimize the risk of false positives given the "proof" gathered and the malicious impact given the intended action would be at step 1.a.3 (see bottom screenshot in Figure 6 above, labeled *How to Navigate to the Protection Results*). This step occurs when a script attempts the first malicious file manipulation. But this information is something that can only be determined after understanding each sub-technique of each stage or test. The success or failure is based on an organization's concern about impeded legitimate user activity vs. the risk of malicious impact from a cyberattack.

## How It Was Stopped Can Make All the Difference in the World

Outside of the discussion of where the attack was stopped, one must consider how the attack was stopped. This is where the analysis takes a tedious turn for the reviewer. For this example, we are going to use three vendor examples of test number two; *TrickBot Execution, Discover, and Kerberoasting*. Out of the 30 vendors that participated in the detection test, 22 participated in the protection test, where test number two wasn't optional. Of the 22 participants, 16 passed the test. We begin to see a better picture of how the 16 stopped that specific attack. To see for yourself, select a vendor from the list, select the 2022 test, and click on protections (see Figure 6). Click on test two and, as mentioned before, look at the section under "Techniques Included in This Test" to see where the attack was stopped. After this, look through the hyperlinked screenshots after clicking on each one so you can zoom in.
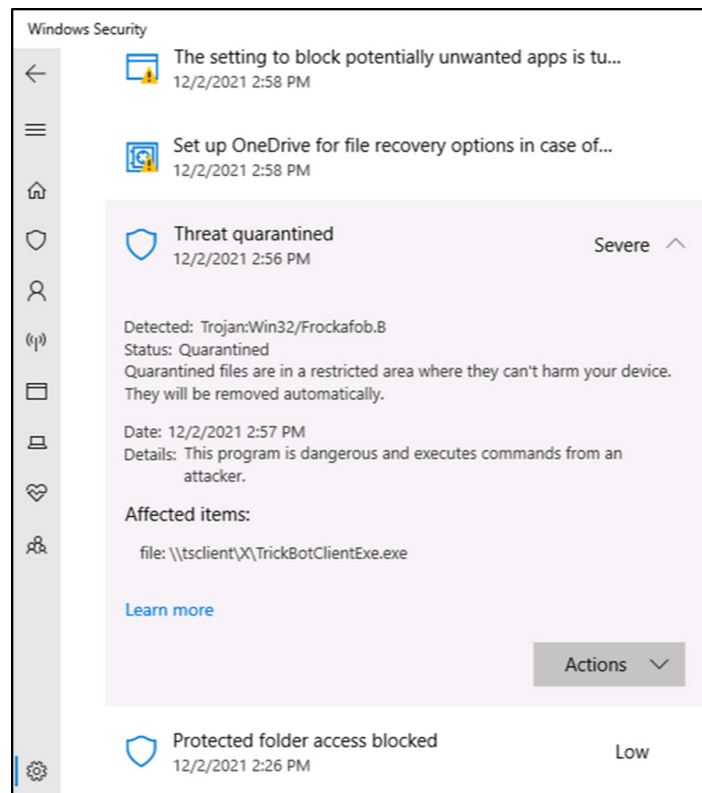


Figure 8: Blocking with Signatures in Test Two

For example, in the case in Figure 8, one vendor blocked the attack on sub-technique 4.A.3 (Ingress Tool Transfer) within test two. When you look at the first screenshot, on the center of the screen, you will see the threat named "Trojan:Win32/ Frockafob.b." When the threat is named, this is your sign that you are dealing with signature-based antivirus (AV), which is only as good as the threat intelligence behind it.  It is also prone to miss new threats until they become known and analyzed, with signatures created/updated. Outside of test two, you will see similar indicators on each of this vendor's  tests (see: Detection source). This approach has been with us from the start of commercial AV applications. While such products are aging out, their insusceptibility to false positives and limited administration is, for some, satisfactory.

Newer than signature-based AV is next-generation AV (NGAV) based on machine learning. You can find many examples of ML to block an attack, which is heavily reliant on the static analysis of malware, among vendors in test two. As one example, the vendor's fourth screenshot for test two shows "Detection type" as "Static." For comparative reasons, you can see similar results in this vendor's first screenshot and note where it talks about the file being written to disk yet convicted by the ML-based AV under "Specific to This Detection." For additional examples of ML detection, check tests 3, 6, and 9 from both vendors as well as 5 and 8 from SentinelOne. Now, this type of detection can also be satisfactory for some organizations. When compared to antivirus, ML-based defenses are stronger at stopping new malware and addressing newer techniques. However, it does requiring regular updating and tuning.
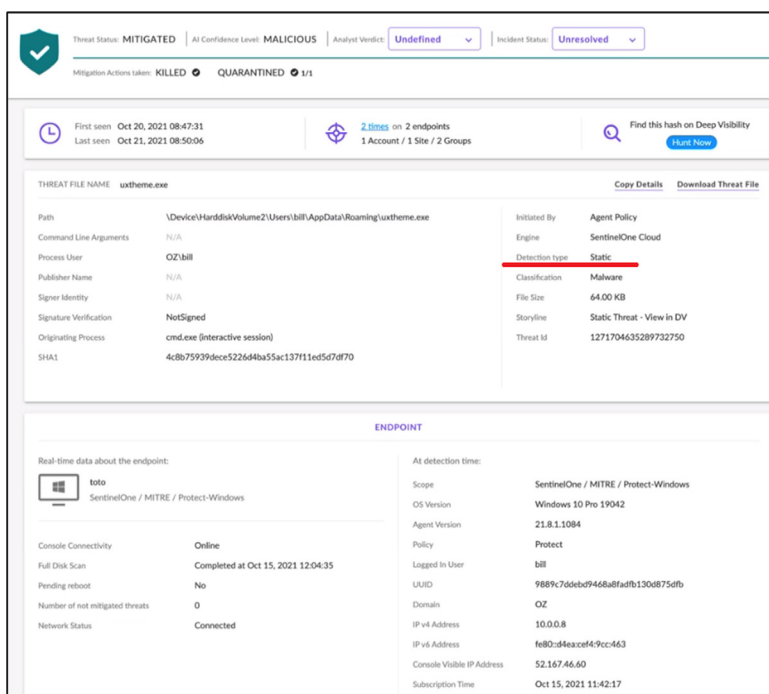


Figure 9: Blocking Malware Using Machine Learning in Test Two

For the strongest protection, against the newest attacks and attack classes, check out an example of behavior-based protections for test two, as shown in the FortiEDR results. In the second screenshot, you can see the result was malicious, along with a verdict from Fortinet Cloud Services (FCS). Fortinet uses FCS on the backend to refine verdicts from the onboard intelligence with FortiEDR. The result is a dynamic or behavior-based approach with zero delays. A behavioral approach is believed to be even stronger when compared to the ML discipline in terms of breadth and depth of prevention. In this specific scenario, the detection is based on the combination of several indicators, the file was copied from remote location, its name is similar to legitimate executable (uxtheme.dll), and it creates suspicious communication immediately after its execution. In terms of test two, FortiEDR stopped the attack at 4.A.4, Application Layer Protocols. Although seeing the attacker manually connect with RDP to transfer a file, this is a benign action used in many applications, and stopping it at this point would produce many false positives. It wasn't until the process arrived at sub-technique number four that FortiEDR realized that the odds of all of these activities all belonging to a benign process are low, and therefore blocked the attack based on behavior.

Furthermore, behavior-based methods like Fortinet solutions offer better protections for zero-day attacks, user-generated malicious installations, and living-off-the-land attacks. Of course, there are drawbacks worth considering. For all the functionality in security they introduce, they do not excel in easing operations and may need to be updated more than the ML approach. This is why many organizations who use a behavior-based EDR solution also rely on their vendor's internal incident response teams to provide a managed EDR (also known as managed detection and response [MDR]) experience to eliminate the burden on their security or security operations center (SOC) staff.
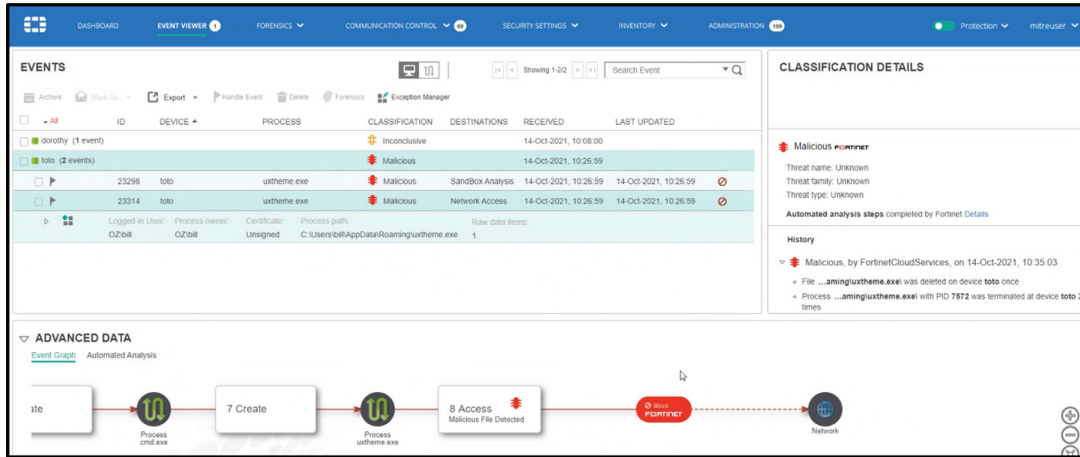


Figure 10: FortiEDR Blocks Using Behavior in Test Two

## Recap on Reading the Report

Reading through the report takes some time and knowledge. This explains why many fall victim to vendor misinformation through easy-to-read charts, stack ranks, and graphics designed to make the report's findings more consumable. The protection tests are simple to diagnose if they were successful or not. Keep in mind that test seven (Linux) was optional. Fortinet chose not to participate because the FortiEDR threat-hunting module for Linux was in beta at the time of the test but was made available with the 5.1 release of the product in Q1 of 2022. Misinformation campaigns will often assign a "Fail" to the five vendors that did not participate in this test when in actuality, eight failed, and eight passed the test. Protection "scores" should assign a percentage of passed attacks out of the number of participating attacks. It is worth noting that if one detected a strong majority of sub-techniques and blocked all attacks on Windows, it doesn't necessarily mean their same success will apply to Linux if they now support that platform. A field test should confirm real-world results. Further, we recommend to go the extra mile, and understand how each vendor you are considering blocked the attack from screenshots. As mentioned, signature and ML static analysis is easier to manage than behavior-based approaches, but more easily bypassed by cybercriminals.

In the detection tests, you can calculate raw detection by dividing the number of total detection minus the number of misses (called "None" on the report) divided by the total number of sub-techniques for the tests they participated in. The total is 109 if they participated in the Linux portion of the test and 90 if they did not. On the overall screen for the 2022 test, The MITRE Foundation created a Visibility rating to make it easier to derive the percentage. Note that the "Detection Count" was a feature in the first three previous tests that was applied to this round.

| Evaluations | Analytic Coverage ⓘ | Telemetry Coverage ⓘ | Visibility ⓘ | Detection Count ⓘ |
|---|---|---|---|---|
| Wizard Spider + Sandworm (2022) | 85 of 90 [2] substeps | 9 of 90 [2] substeps | 87 of 90 [2] substeps | ⊖ |

Figure 11: FortiEDR Overall Detection Results for the Wizard Spider + Sandworm MITRE ATT&CK Evaluation

Once the overall detection score is understood, you should create a percentage for coverage. In the example above, that would be 85/90 as the Linux test, with its 19 additional sub-techniques, was not in scope. This rounds to 94%, but misinformation campaigns can sometimes divide 85 by 109 to give the appearance that this vendor found 78% of the right techniques used in the evaluation.

It is also worth pointing out that those with stronger telemetry coverage results are arguing that this is a "purer style of detection," but can be argued that The MITRE Foundation would not agree with that assessment. Solutions that produce a lot of telemetry are typically not good for an organization that doesn't have a security team or one that has a time-constrained security or SOC team, as it requires much more work to diagnose a threat.

## What Style of Endpoint Security Is Right for Your Organization?

The approach you take to security depends on what is right for your organization. As mentioned above in the section on How It Can Be Stopped, there are different approaches that your organization must consider. When comparing the top EPP solutions together, one can't only go by the flat protection, visibility, and analytic scores within the MITRE test, especially when the percentages are all within a few points of each other. We recommend reading the Gartner® report "Comparison of the Impacts of Endpoint Protection Techniques"[1] to better understand the benefits and drawbacks to each type of security strategy.
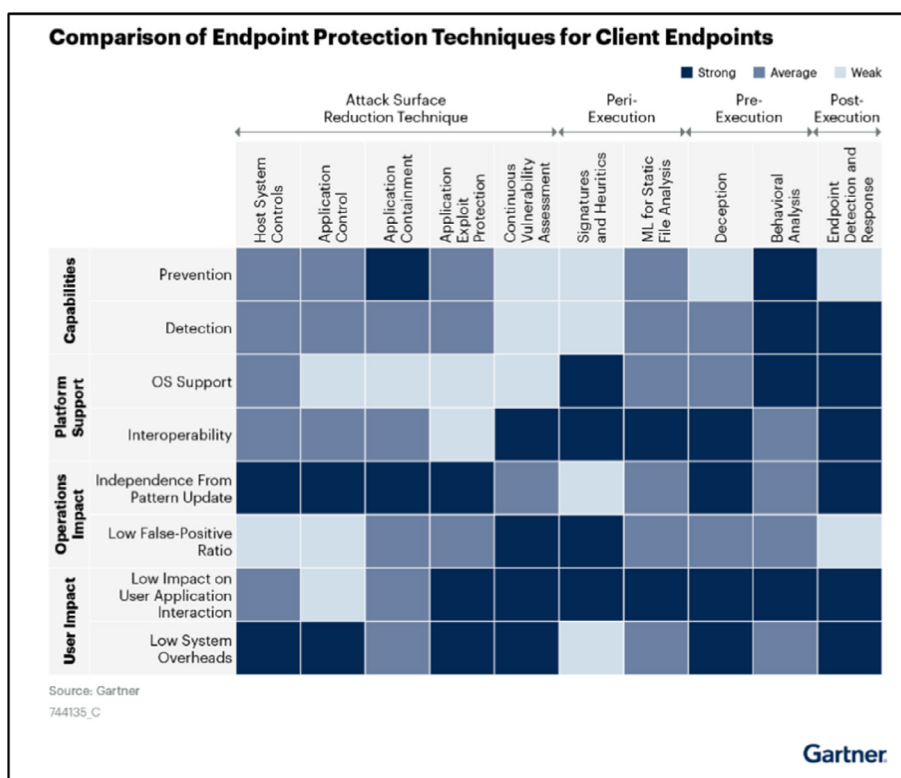


Figure 12: Gartner® comparison of Endpoint Protection Techniques for Client Endpoints

## How Did Fortinet Do?

In the fourth round of the MITRE ATT&CK Enterprise Evaluations, Fortinet showed how much it had improved from the year before by redesigning the client to better speak along the lines of the MITRE framework. For the second year in a row, **FortiEDR blocked all attacks**. It is worth noting again that because the Linux threat-hunting portion was in beta at the time of the test (but now available), it did not participate in test seven. Upon evaluating the screenshots provided in the test, you can see that FortiEDR blocked six of eight participating tests with a behavior-based approach and two with an ML approach. Regardless of technique, all attacks were blocked without signatures and relied on the onboard behavior-based intelligence that makes it excel in head-to-head comparisons, especially when offline. Read this data sheet for more information on FortiEDR.

When it came to detections, the visibility score was 97% (87/90), which was one of the top-five results in the test in terms of percentage. In addition to that, 94% (85/90) of all sub-techniques were detected with analytics, which was also a top-five result from the test also based on percentage. Analytics shows the ability to block unknown and zero-day attacks, which is why this metric is so highly regarded and examined by students of the report.

We look forward to the fifth round of the MITRE Evaluations, where FortiEDR will participate in any Linux-based tests. The threat-hunting feature has been made available and matured for the Linux environment. If you like researching the performance of EPP solutions in real-world scenarios, review the non-sponsored research from the University of Piraeus titled "An Empirical Assessment of Endpoint Security Systems Against Advanced Persistent Threats Attack Vectors." This 57-page paper by two IT security experts is in its third and final version. They attempted to bypass the world's EDR solutions to prove or disprove their efficacy in four real-world sets of attacks. It is worth noting that FortiEDR was the first EDR solution out of the box to block all of their attacks in their second round and, by the third, was only one of two.

| EDR | CPL | HTA | EXE | DLL |
|---|---|---|---|---|
| BitDefender GravityZone Plus | ✗ | ✗ | ✓ | ✗ |
| Carbon Black Cloud | ⋆ | ⋆ | ✓ | ✓ |
| Carbon Black Response | ● | ✗ | ✓ | ✓ |
| Check Point Harmony | ✗ | ◇ | ✗ | ✓ |
| Cisco AMP | ✗ | ✗ | ✓ | ⊙ |
| Comodo OpenEDR | ✗ | ✓ | ✗ | ✓ |
| CrowdStrike Falcon | ✓ | ✓ | ✗ | ✓ |
| Cylance PROTECT | ○ | ○ | ✓ | ✗ |
| Cynet | ✗ | ✓ | ✓ | ✓ |
| Elastic EDR | ✗ | ✓ | ✓ | ✗ |
| F-Secure Elements Endpoint Detection and Response | ◇ | † | ✓ | ✗ |
| FortiEDR | ✗ | ✗ | ✗ | ✗ |
| Harfang Lab Hurukai | ✗ | ✓ | ✗ | ✓ |
| ITrust ACSIA | ✓ | ✓ | ✓ | ✓ |
| McAfee Endpoint Protection with MVision EDR | ✗ | ● | ✓ | ✓ |
| Microsoft Defender for Endpoints (original IOCs) | ⋆ | ✗ | ✗ | ✓ |
| Microsoft Defender for Endpoints (Updated MDE) | ⋆ | ✗ | ✗ | ✗ |
| Microsoft Defender for Endpoints (Updated MDE & IOCs) | ∇ | ✗ | ✗ | ✓ |
| Minerva Labs | ⊕ | ✗ | ✓ | ✗ |
| Palo Alto Cortex | ✓ | ✓ | ✗ | ✓ |
| Panda Adaptive Defense 360 | ✗ | ✓ | ⋆ | ✓ |
| Sentinel One (Original version) | ✓ | ✓ | ✓ | ✗ |
| Sentinel One (Current Version) | ✗ | ✗ | ✗ | ✗ |
| Sophos Intercept X with EDR | ✗ | ✗ | ✓ | - |
| Symantec Endpoint Protection Complete | ⋆ | ✗ | ⋆ | ⋆ |
| Trend micro Apex One | ● | ● | ✓ | ✓ |
| **Endpoint Protection** | | | | |
| ESET PROTECT Enterprise | ✗ | ✗ | ✓ | ✓ |
| F-Secure Elements Endpoint Protection Platform | ✓ | ✓ | ✓ | ✓ |
| Kaspersky Endpoint Security | ✗ | ✗ | ✗ | ✓ |
| McAfee Endpoint Protection | ✗ | ✗ | ✓ | ✓ |
| Symantec Endpoint Protection | ✓ | ✗ | ✓ | ✓ |

Figure 13: X Marks the Spot of a Failed Attack on Page 44 of the University of Piraeus Research Report[2]

[1] Shashank Sharma and Mario de Boer, "Comparison of the Impacts of Endpoint Protection Techniques," Gartner, June 30, 2021.

[2] George Karantzas and Constantinos Patsakis, "An Empirical Assessment of Endpoint Security Systems Against Advanced Persistent Threats Attack Vectors," Department of Informatics, University of Piraeus, Greece, July 9, 2021.

GARTNER is a registered trademark and service mark of Gartner, Inc. and/or its affiliates in the U.S. and internationally and is used herein with permission. All rights reserved.

**F⊟RTINET**®

www.fortinet.com