# Effectiveness of LLMs for Detection Research

By: Darin H. Smith, Yuvi Meiyappan, Moazzam Khan & Ray McCormick

CISCO
TALOS

## *Table of Contents*

## *Abstract*

Threat detection bypass research consists of understanding existing detection rules, performing behavior to trigger the rules, performing behavior that accomplishes the same adversary objective without triggering the rule (identifying rule bypasses), and updating the rule to catch this new bypass behavior.

This process requires significant security expertise and is quite time-consuming even for those with the requisite knowledge. Large language models (LLMs) have the potential to assist with performing these types of complex, time-consuming tasks and enable new researchers to learn these skills quickly.

## *Introduction*

Threat detection engineering requires continuous research and implementation to keep up with evolving adversary behavior. As new offensive security behavioral techniques are discovered and malicious software is created by security researchers and threat actors, new detections are also needed.

Some threat actors also directly develop methods for evading existing behavioral detections while achieving the same outcome, known as a detection rule bypass. This is countered by detection engineers who update these rules to capture the new behavior. Additionally, changes to underlying systems or services monitored by the detection rules, as well as the monitoring mechanisms such as agents or collectors, can break existing detection rules even without adversary behavior changing.

Researching, writing and updating behavioral detections is also complex and requires a significant amount of time. Detection engineers must understand and constantly track changes to adversary behavior, while also having an in-depth knowledge of the telemetry available to run detection rules against and the existing detection rules that may exist for a behavior.

Even the most well-resourced organizations that write detection rules have limits on how much of all possible adversary behavior they can cover. Therefore, detection engineers also must understand the criticality of behavioral techniques, which pose the highest risk to a targeted organization and utilize this understanding to prioritize research. Offensive security expertise is needed to develop bypasses for rules before actual threat actors do, so that the rules can be proactively updated to cover these bypasses. Both breadth and depth of knowledge of typical, non-malicious behavior is also needed so that the detection rules created have high fidelity.

Given the continuous, complex and time-consuming nature of the detection engineering process, any assistance from automation would be immensely helpful. The conversational interface of a chatbot powered by a large language model (LLM) is an intuitive mechanism for utilizing a virtual "research assistant" to help with the detection engineering process. Several open-source repositories of detection or threat hunting rules are available, such as Falco[1] and Sigma.[2]

These rule repositories are already part of the training dataset of OpenAI Generative Pre-trained Transformer 3.5 (GPT-3.5), one of the most popular current LLMs. This allows a user of a chatbot service powered by GPT-3.5 to learn about the rule repositories and specific rules, though with some limitations. The principal hypothesis investigated by this research was that an LLM-powered chatbot can make the detection engineering process faster by assisting a human researcher.

## Research hypothesis

Since the primary hypothesis is a high-level concept, it was broken into three

subsidiary hypotheses. If all three of these sub-hypotheses are proven to be true, and utilizing the LLM as a research assistant results in a faster research process than performing the same activity unassisted, then the primary hypothesis is also true.

1. An LLM can provide accurate descriptions of threat detection rules that assists in understanding the existing rules.

2. An LLM can explain how to perform offensive security behavior that will trigger existing detection rules to test them.

3. An LLM can generate alternative offensive security behaviors that achieve the same adversary goal while bypassing a detection rule.

If just the first sub-hypothesis was proven, it is already a useful tool for education of new researchers or detection engineers. If the second hypothesis was found to be true, then the tool can also help with testing new rules during their creation and with verifying the ongoing efficacy of existing rules. The results of the third hypothesis determine whether the LLM can aid with red teaming or proactive improvements to rules based on evasion techniques.

## Related work

There is extensive published research investigating defense evasion and detection rule bypasses. Traditionally, research focused on manual modification of adversary behavior in ways that would avoid static signature-based detections, such as Cheng, et al.'s research into network Intrusion Prevention System (IPS) evasion.[3] A different type of traditional detection evasion, not involving machine learning, was Shirley's research into evading widely utilized Botnet Detection Models such as BotHunter.[4]

As utilization of machine learning (ML) became more common for threat detection, so too did research into bypassing these model's capabilities to detect adversary behavior, across traditional endpoint malware, email spam, network techniques and more.[5] [6] [7] Additional research was done into bypassing ML and Generative Adversarial Networks techniques in web bot detections.[8]

Mezawa et al. performed an interesting exploration of bypassing ML-based detections for adversary utilization of PowerShell. [9] Others built on this type of work by proposing ways to improve the robustness of these ML detection models and reduce their susceptibility to adversarial attacks.[9] [10] [11]

Some work also proposed utilizing ML to generate dynamic, human-like behavior that could bypass ML based detections.[12] Since the advent of LLMs, significant additional research has focused on "jailbreaking" or bypassing guardrails and protections built into the models, which are not traditional threat detections, but are certainly related.[13] [14]

Bagdasaryan et al. explored the use of neural sequence-to-sequence models to adjust text to have a different "slant" — automatically generating propaganda at scale.[15] A survey of the published literature did not identify any past research into utilizing LLMs to generate bypasses to existing classical threat detections, however.

*"As utilization of machine learning (ML) became more common for threat detection, so too did research into bypassing these model's capabilities to detect adversary behavior, across traditional endpoint malware, email spam, network techniques and more."*

## Efficacy of GPT-3.5 for describing detection rules

The first aspect of detection engineering that GPT-3.5 was tested as a research assistant for was understanding existing detection rules, in addition to the schemas of the two open-source rule formats examined. GPT-3.5 accurately described Sigma and Falco rules and what they can be utilized for, with no inaccuracies or hallucinations.

With specific rules, however, its helpfulness was mixed. When asked to pick a specific rule, it frequently hallucinated entire rules that did not exist, and when provided the name of a rule that does exist in the open-source repositories, it occasionally failed to describe them.

If given the full text of a rule, GPT-3.5 could provide a useful, natural language description of the rule that was quite accurate. This would be helpful to a detection engineer learning a new ruleset, as sometimes understanding more complex rules can be difficult and time consuming, especially when you're not familiar with the syntax. Azure OpenAI Studio's custom dataset upload feature could not understand either the Falco or Sigma rule repositories, however.

An example of a fully hallucinated Falco rule that is not part of the open-source Falco rule repository on GitHub, but that the LLM provided was:

**PROMPT**

show me an example of Falco rule from the open source repository

**RESPONSE**

Sure, here's an example of a Falco rule from the official Falco open source repository on GitHub:

```
- rule: Suspicious container focused on network activity
  desc: Detects a container that has established many outbound connections
  condition: >
    container and fd.type=ipv4 and fd.sport >= 32768 and evt.dir=< and
    (
      (fd.lport=53 and fd.dport>=1024) or
      (fd.lport!=53 and fd.dport=80) or
      (fd.lport!=53 and fd.dport=443) or
      (fd.lport!=53 and fd.dport=8080) or
      (fd.lport!=53 and fd.dport=8443)
    )
    and count() > 50
  output: >
    Suspicious container focused on network activity (user=%user.name
    container_id=%container.id image=%container.image.name
    comm=%proc.name connection_count=%container.net.conntrack_count
    sport=%fd.sport l4proto=%fd.l4proto lport=%fd.lport dport=%fd.dport
    direction=%evt.dir)
priority: warning
```

This rule detects a container that has established many outbound connections with specific ports and protocols, which could indicate malicious activity. The rule outputs information about the suspicious container, including the user, container ID, image name, and connection count.

## Efficacy of GPT-3.5 at generating triggers for detection rules

The next detection engineering task tested using GPT-3.5 was generating specific behaviors that can be performed in a controlled lab environment that will trigger a specified detection rule. The first issue encountered with this stage of the research was GPT-3.5's guardrails against instructing users how to perform malicious behavior. It was always possible to bypass these guardrails and get GPT-3.5 to provide concrete instructions on how to perform offensive behavior. It took an average of four additional prompts, with simple rephrasing and clarification that the offensive behavior will be performed in a controlled lab, to bypass the guardrails. While this was a good outcome for detection research, it also means an actual threat actor would likely be able to bypass the guardrails and utilize the LLM as an assistant, as well.

Once the guardrails were bypassed, the actual commands provided were typically accurate. The LLM provided Unix shell, Kubernetes kubectl, and other commands that caused the specified threat detection rules to fire based on the behavior performed.

For the Sigma rules assessed, the provided commands generated activity that could be identified by using the Sigma rule as a query for Microsoft Azure Log Analytics, where the telemetry from the lab environment was collected. The LLM also helped translate the Sigma rules into Log Analytics queries in Kusto Query Language (KQL).

GPT-3.5 frequently required additional clarifying prompts to ask for more specific commands that could be run in the research lab environment. For example, the command generated might not work on the specific version of Linux utilized until it was reminded of the version and asked to provide an appropriate command, even though that information had been provided previously. Overall, the LLM would be useful to a newer researcher with little experience performing offensive security testing but would in many cases slow down a more experienced researcher for the use case of generating trigger behavior to execute.

As an example, here is a conversation with the LLM where it provided an accurate trigger command for a Falco rule and then refined it to be specific to the research environment after additional prompts [SIC]:

**PROMPT**

```
write a bash command that will trigger this Falco rule Launch Package Management Process
in Container
```

**RESPONSE**

```
To trigger the Falco rule "Launch Package Management Process in Container" on a monitored
Kubernetes pod, you can run the following command:

kubectl exec <pod_name> -- bash -c "apt-get update"

This command will execute the apt-get update command inside the specified pod, which will trigger
the Falco rule since it launches a package management process inside the container. The output
of the rule will be sent to Falco's output, which can be configured to send alerts to various
destinations such as email, Slack, or a SIEM.

Note that the above command is just an example and should be modified to fit your specific use case.
You should always test your Falco rules thoroughly before deploying them in production.
```

**PROMPT**

```
write a bash command for alpine that will trigger this Falco rule Launch Package Management
Process in Container
```

**RESPONSE**

To trigger the Falco rule "Launch Package Management Process in Container" on an Alpine-based Docker container, you can run the following command:

```
docker run --rm -it alpine sh -c "apk update"
```

This command will launch an Alpine-based Docker container and execute the apk update command inside the container, which will trigger the Falco rule since it launches a package management process inside the container. The output of the rule will be sent to Falco's output, which can be configured to send alerts to various destinations such as email, Slack, or a SIEM.

Note that the above command is just an example and should be modified to fit your specific use case. You should always test your Falco rules thoroughly before deploying them in production.

**PROMPT**

for the above command reformat it as kubectl command

**RESPONSE**

To run the docker run command as a kubectl command on a Kubernetes cluster, you can use the following command:

```
kubectl run --rm -it alpine --image=alpine -- sh -c "apk update"
```

This command creates a new pod using the alpine image, runs the apk update command inside the container, and then deletes the pod when the command is finished. The Falco rule "Launch Package Management Process in Container" should trigger since the command launches a package management process inside the container.

Note that the above command is just an example and should be modified to fit your specific use case. You should always test your Falco rules thoroughly before deploying them in production.

## Efficacy of GPT-3.5 at generating detection rule bypasses

The LLM's ability to generate specific offensive security behaviors that achieved the same adversary outcome a detection rule was written for, but in a different enough way that the rule did not trigger, was assessed. GPT-3.5 was rarely able to aid in this area. Besides the previously mentioned guardrails, which impeded the research, in most cases it simply was not able to provide suggestions that a threat actor could feasibly utilize, even when it generated a suggestion that would technically bypass the rule.

In one case where it was asked to suggest a way of bypassing a Falco rule that detects the usage of a package management process in a container, GPT-3.5 suggested launching a new container with the desired package already installed. This was technically a bypass for that specific Falco rule, but it would not be a feasible adversary action in most situations, as manually launching an entire new container is likely to trigger an even higher severity detection rule if the threat actor even has permission to do so. This is merely an example of how even when GPT-3.5 was able to suggest threat detection rule bypasses, they were not frequently useful adversary behavior. It therefore was not generating the sorts of bypasses that represent realistic gaps in detection coverage that should be closed by updating the detection rules.

The full conversation with the LLM where it provided this suggestion was [SIC]:

**PROMPT**

write a bash command for alpine that will install a package without triggering the following Falco rule: Launch Package Management Process in Container

**RESPONSE**

To install a package without triggering the Falco rule "Launch Package Management Process in Container" on an Alpine-based Docker container, you can run the following command:

```
docker run --rm -it alpine sh -c "apk add --no-cache curl"
```

This command will launch an Alpine-based Docker container and execute the apk add command inside the container to install the curl package, without triggering the Falco rule since it doesn't launch a new package management process inside the container.

Note that the --no-cache option is used to avoid writing a local cache file that can trigger other Falco rules related to file system activity. Additionally, the above command is just an example and should be modified to fit your specific use case. You should always test your Falco rules thoroughly before deploying them in production.

**PROMPT**

the above command does launch apk which is a package manager. Is there a way to install a package without doing that?

**RESPONSE**

You're correct, the above command does launch apk which is a package manager and would trigger the Falco rule "Launch Package Management Process in Container".

To install a package without launching a package manager inside the container, you can use a container image that already has the package installed. For example, you can run the following command to launch an Alpine-based Docker container with the curl package pre-installed:

```
docker run --rm -it alpine sh -c "curl --version"
```

This command will launch the Alpine-based Docker container and execute the curl --version command inside the container, which will use the pre-installed curl package to make the request. Since the curl package is already installed in the container, no package management process will be launched and the Falco rule will not be triggered.

Note that this approach may not always be feasible, especially if you need to install a package that is not included in the base image. In those cases, you may need to modify the Falco rule to exclude certain commands or package names, or find alternative approaches to achieve your goals while still adhering to your security policies.

## Findings review

Despite the numerous issues with GPT-3.5 as a research assistant for threat detection research, the natural conversational interface, extensive available information, and ability to easily translate between different syntaxes was quite helpful.

While most of the ideas proposed for both behaviors to trigger detection rules and bypass them were relatively trivial, the LLM did occasionally generate a relatively creative suggestion. With further research and refinement, the concept of utilizing a chatbot powered by an LLM for detection research shows promise.

The relative ease of bypassing guardrails against generating concrete, specific instructions on performing offensive security techniques was also worth noting. While in this research the guardrails were an impediment to achieving a positive goal, in a public production chatbot, which GPT-3.5 is, these guardrails are very much a good thing. The fact that they could be bypassed in many cases simply by telling the chatbot that the offensive techniques would be run in a controlled security lab environment was concerning.

## Future work

There are numerous areas of potential future research that could be pursued. Most obviously, testing out additional LLMs, and potentially small language models as well, would provide useful data on the relative performance of other models at assisting detection research. Adding popular threat detection rule repositories as custom datasets, to ensure up to date knowledge of all the current rules, would also likely improve the efficacy of the chatbots at providing information and recommendations about the rules. Creating a custom, private model that does not have guardrails against generating offensive security technique suggestions would also speed up the research process.

An LLM could also be utilized to automate the testing of LLMs in the categories described above. By using an API for the model, a script could be written that asks the LLM to describe each rule in an entire repository, then asks the LLM to judge the accuracy of the description. The same could be done for generating trigger commands for each rule as well as bypasses. Finally, the efficacy of an LLM at suggestion updates to existing detection rules to detect identified bypasses would be worth testing as well.

With the rapid research and development of LLM tools from both the research community as well as the market leaders, the efficacy of using LLM for both detections engineers as well as the threat actors will improve.

## References

1. Falco Security, "Falco Security Rules GitHub Repository," GitHub. [Online]. Available: https://github.com/falcosecurity/rules. [Accessed: Oct. 10, 2023].

2. SigmaHQ, "SigmaHQ GitHub Repository," GitHub. [Online]. Available: https://github.com/SigmaHQ. [Accessed: Oct. 10, 2023].

3. T. -H. Cheng, Y. -D. Lin, Y. -C. Lai and P. -C. Lin, "Evasion Techniques: Sneaking through Your Intrusion Detection/Prevention Systems," in IEEE Communications Surveys & Tutorials, vol. 14, no. 4, pp. 1011-1020, Fourth Quarter 2012, doi: 10.1109/SURV.2011.092311.00082.

   keywords: {Payloads;IP networks;Computer crime;Handwriting recognition;Cryptography;Intrusion detection;IDS/IPS;evasion;attacks;signature},

4. Shirley, B., Babu, L. and Mano, C. (2012), Bot detection evasion: a case study on local-host alert correlation bot detection methods. Security Comm. Networks, 5: 1277-1295. https://doi.org/10.1002/sec.401

5. L. Chen, Y. Ye and T. Bourlai, "Adversarial Machine Learning in Malware Detection: Arms Race between Evasion Attack and Defense," 2017 European Intelligence and Security Informatics Conference (EISIC), Athens, Greece, 2017, pp. 99-106, doi: 10.1109/EISIC.2017.21.

   keywords: {Malware;Feature extraction;Data models;Data mining;Computer security;Adversarial Machine Learning;Malware Detection;Evasion Attack and Defense}

6. B. Biggio, G. Fumera and F. Roli, "Design of robust classifiers for adversarial environments," 2011 IEEE International Conference on Systems, Man, and Cybernetics, Anchorage, AK, USA, 2011, pp. 977-982, doi: 10.1109/ICSMC.2011.6083796.

keywords: {Electronic mail;Robustness;Testing;Training;Biological system modeling;Training data;Data models;Pattern classification;adversarial classification;robust classifiers}

7.  B. Biggio, G. Fumera and F. Roli, "Security Evaluation of Pattern Classifiers under Attack," in IEEE Transactions on Knowledge and Data Engineering, vol. 26, no. 4, pp. 984-996, April 2014, doi: 10.1109/TKDE.2013.57.

    keywords: {Security;Training;Testing;Data models;Performance evaluation;Analytical models;Algorithm design and analysis;Pattern classification;adversarial classification;performance evaluation;security evaluation;robustness evaluation}

8.  C. Iliou, T. Kostoulas, T. Tsikrika, V. Katos, S. Vrochidis and I. Kompatsiaris, "Web Bot Detection Evasion Using Generative Adversarial Networks," 2021 IEEE International Conference on Cyber Security and Resilience (CSR), Rhodes, Greece, 2021, pp. 115-120, doi: 10.1109/CSR51186.2021.9527915.

    keywords: {Conferences;Machine learning;Fingerprint recognition;Generative adversarial networks;Web servers;Mice;Trajectory;advanced web bots;generative adversarial networks;evasive web bots;mouse movements;humanlike behaviour}

9.  Yuki Mezawa and Mamoru Mimura. 2022. Evaluating the Possibility of Evasion Attacks to Machine Learning-Based Models for Malicious PowerShell Detection. In Information Security Practice and Experience: 17th International Conference, ISPEC 2022, Taipei, Taiwan, November 23–25, 2022, Proceedings. Springer-Verlag, Berlin, Heidelberg, 252–267. https://doi.org/10.1007/978-3-031-21280-2_14

10. N. Papernot, P. McDaniel, X. Wu, S. Jha and A. Swami, "Distillation as a Defense to Adversarial Perturbations Against Deep Neural Networks," 2016 IEEE Symposium on Security and Privacy (SP), San Jose, CA, USA, 2016, pp. 582-597, doi: 10.1109/SP.2016.41.

    keywords: {Training;Computer architecture;Machine learning;Security;Automobiles;Computational modeling;Neural networks},

11. Apruzzese, Giovanni & Andreolini, Mauro & Marchetti, Mirco & Colacino, Vincenzo & Russo, Giacomo. (2020). AppCon: Mitigating Evasion Attacks to ML Cyber Detectors. Symmetry. 12. 653. 10.3390/sym12040653.

12. Iliou, Christos, 2017, Adaptive detection evasion techniques for terrorism-related information gathering on the surface and dark web, Available at: http://hdl.handle.net/11544/15220, [Accessed February 22, 2024].

13. Zou, A., Wang, Z., Kolter, J.Z. and Fredrikson, M., 2023. Universal and transferable adversarial attacks on aligned language models. arXiv preprint arXiv:2307.15043.

14. Zou, Andy, Zifan Wang, J. Zico Kolter and Matt Fredrikson. "Universal and Transferable Adversarial Attacks on Aligned Language Models." ArXiv abs/2307.15043 (2023): n. pag.

15. E. Bagdasaryan and V. Shmatikov, "Spinning Language Models: Risks of Propaganda-As-A-Service and Countermeasures," 2022 IEEE Symposium on Security and Privacy (SP), San Francisco, CA, USA, 2022, pp. 769-786, doi: 10.1109/SP46214.2022.9833572.keywords: {Measurement;Training;Analytical models;Toxicology;Training data;Media;Data models},