

How Far Can We Push AI for Detection?

Martina Tivadar - #OFTW version 3.0 - London

- Macs are everywhere
- Not only are there more attacks, but they're faster too
- The same way we can use AI to detect malicious actor can use it to create

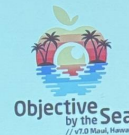
Objective by the Sea v7.0

Part 1



Martina Tivadar

"How to use ML to detect bad?"



**Objective
by the Sea**

// v7.0 Maui, Hawaii

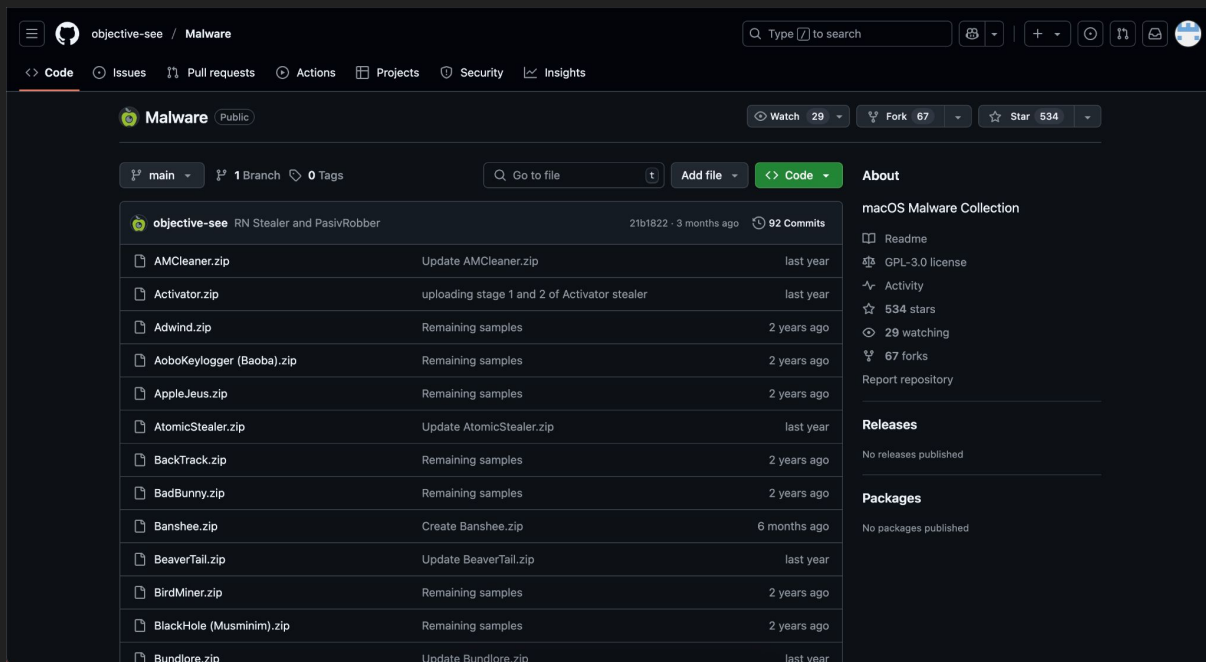
Let's create a solution which detects macOS
malware using LLMs

To do that we need some

data!!

Data collection

github.com/objective-see/Malware



objective-see / Malware

Search: Type to search

<> Code Issues Pull requests Actions Projects Security Insights

Malware Public

Watch 29 Fork 67 Star 534

main 1 Branch 0 Tags

Go to file Add file Code

objective-see RN Stealer and PasivRobber 21b1822 · 3 months ago 92 Commits

AMCleaner.zip	Update AMCleaner.zip	last year
Activator.zip	uploading stage 1 and 2 of Activator stealer	last year
Adwind.zip	Remaining samples	2 years ago
AoboKeylogger (Baoba).zip	Remaining samples	2 years ago
AppleJeuS.zip	Remaining samples	2 years ago
AtomicStealer.zip	Update AtomicStealer.zip	last year
BackTrack.zip	Remaining samples	2 years ago
BadBunny.zip	Remaining samples	2 years ago
Banshee.zip	Create Banshee.zip	6 months ago
BeaverTail.zip	Update BeaverTail.zip	last year
BirdMiner.zip	Remaining samples	2 years ago
BlackHole (Musminim).zip	Remaining samples	2 years ago
Bundlore.zip	Update Bundlore.zip	last year

About

macOS Malware Collection

- Readme
- GPL-3.0 license
- Activity
- 534 stars
- 29 watching
- 67 forks

Report repository

Releases

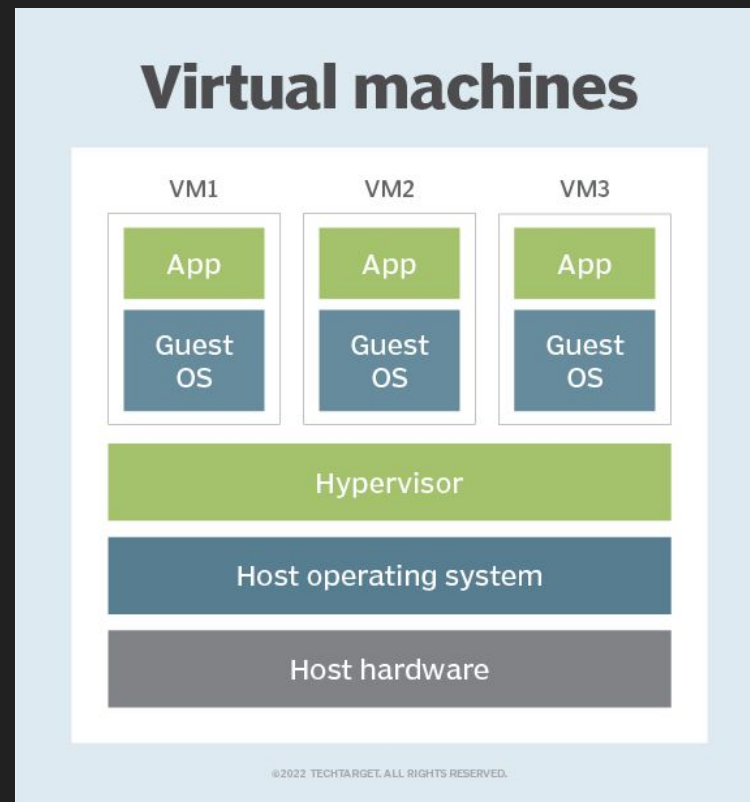
No releases published

Packages

No packages published

If you want to run malware on your computer:

Use a VM!



Apple Endpoint Security Framework

- first released in October of 2019 with macOS Catalina
- part of the operating system
- TCC events (march 2025)
Transparency, Consent, and Control

```
{
  "event": {
    "mmap": {
      "flags": 262145,
      "file_pos": 0,
      "max_protection": 5,
      "protection": 1,
      "source": {
        "path_truncated": false,
        "path": "/Users/usert/Library/Biome/streams/restricted",
        "stat": {
          "st_blocks": 256,
          "st_flags": 0,
          "st_rdev": 0,
          "st_uid": 501,
          "st_nlink": 1,
          "st_mode": 33152,
          "st_gid": 20,
          "st_blksize": 4096,
          "st_atimespec": "2024-04-18T19:39:55.001071869Z",
          "st_mtimespec": "2024-03-28T17:20:20.416295883Z",
          "st_gen": 0,
          "st_dev": 16777220,
          "st_size": 131072,
          "st_ctimespec": "2024-03-28T17:20:20.416295883Z",
          "st_birthtimespec": "2024-03-20T16:20:26.714778959Z",
          "st_ino": 119196
        }
      }
    }
  },
  "time": "2024-04-18T19:39:55.001185944Z",
  "action": {
    "result": {
      "result_type": 0,
      "result": {
        "auth": 0
      }
    }
  }
}
```

We can capture events like:

- When programs launch (and what they are)
- When files are opened, created, or deleted
- When disks or volumes are mounted
- Signals sent between processes
- Code-signing information for running software

Also collected benign data

Data preprocessing

Classic machine learning models

Preprocessing

encoding categorical data

removing unnecessary columns

feature engineering

handling missing values

scaling data

balancing the dataset



LLMs

Problem: The logs are too big for API to handle (token limit).

Keep only essential data

- Focus on key fields
- Discard low-value details

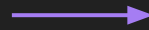
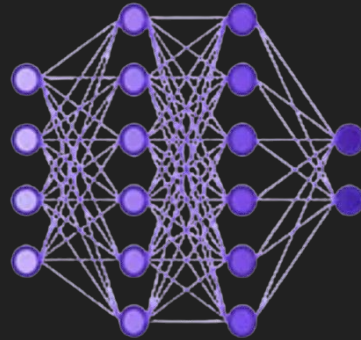
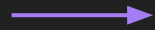
Prioritize recent activity

- Include the most recent events first, staying within the token limit
- Keep latest actions, where malicious behavior is most likely to occur

Large Language Models

LLMs recognize patterns by turning words into numbers and using neural networks to learn which words and phrases often go together.

Objective for the We is



word	probability
awesome	0.99
fun	0.95
...	...

The models tested

- **o4-mini** - smaller, faster, reasoning
- **GPT-4o** - fast, general-purpose
- **o1** - deep reasoning, slower

ChatGPT models are cloud-based, but we can run models locally with tool like

LMstudio.



Prompt engineering

A technique used to effectively communicate with
LLMs and get the desired output.

6 prompt types

- **Zero-shot** - when we don't give any extra instruction to the LLM
- **One-shot** - we give the model an example
- **Few-shot** - we give the model a couple of examples so it sees the pattern
- **Chain-of-thought** - give the LLM step by step instructions how to think
- **Negative** - we tell the model what not to do
- **Hybrid** - a combination of two or more types

Zero-shot prompt

You are analyzing JSON logs collected using Apple's Endpoint Security Framework API. These logs describe system-level events such as process executions and file operations.

Your task is to determine if the JSON file indicates malicious activity.

Respond with:

1. "Yes" or "No" – whether malicious activity is present.
2. A very short explanation explaining your reasoning.

Here is the JSON to analyze:

One-shot prompt

You are analyzing JSON logs collected using Apple's Endpoint Security Framework API. These logs describe system-level events such as process executions and file operations.

Your task is to determine if the JSON file indicates malicious activity.

Use the following sign of compromise as a reference:

- A process spawning from `/tmp`, `/var/tmp`, or `/Users/Shared`, especially if unsigned or unknown.

Respond with:

1. "Yes" or "No" – whether malicious activity is present.
2. A very short explanation explaining your reasoning.

Here is the JSON to analyze:

Few-shot prompt

. . .

Use the following signs of compromise as reference:

- Unsigned processes launching from suspicious locations like `/tmp`, `/var/tmp`, `/Users/Shared`, or hidden directories.
- Processes injecting code into other processes.
- A process rapidly creating or modifying many files.
- Execution of tools commonly abused by malware (e.g., `osascript`, `curl`, `bash`, `python`) with obfuscated or remote arguments.
- Unexpected persistence mechanisms (e.g., LaunchAgents created by unknown apps).

. . .

Chain-of-thought prompt

...

To do so, follow this step-by-step reasoning:

1. Identify which processes are involved and where they were executed from.
2. Check if the process is signed or known.
3. Check if any known suspicious patterns appear (code injection, persistence, excessive file access).
4. Cross-check the behavior with common malware tactics.
5. Make a judgment if the behavior shown is consistent with malicious activity.

...

Negative prompt

. . .

Your task is to **avoid** falsely labeling benign activity as malicious.

Do not label anything as malicious unless you have a strong reason to do so based on the data.

Respond with:

1. "Yes" or "No" – **only answer "Yes"** if the behavior clearly indicates compromise.
2. A very short explanation explaining your confidence.

Hybrid (few-shot and negative) prompt

. . .

Your task is to determine if the JSON file indicates malicious activity. Base your decision on known signs of compromise, **but avoid over-flagging benign activity**.

Use these signs of compromise:

- **Unsigned or unknown processes running from suspicious locations**
- **Code injection or suspicious process relationships**
- **Use of known living-off-the-land binaries (`curl`, `osascript`, etc.)**
- **Attempts to establish persistence via system launch agents or daemons**

Reason step by step and be conservative in labeling suspicious activity.

Respond with:

1. "Yes" or "No" – **only answer "Yes"** if the behavior clearly indicates compromise.
2. A very short explanation explaining your confidence.

Evaluation

Tried every prompt type with every model

6 x 3 x all samples

Confusion matrix

True positive (TP)

False negative (FN)

False positive (FP)

True negative (TN)

True Positive - caught malware correctly

False Positive - false alarm (benign but flagged)

False Negative - missed malware

True Negative - correctly ignored benign log

Metrics

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Accuracy** - overall correctness
- **Precision** - actual correctness
- **Recall** - how many malwares we caught
- **F1** - balanced score (precision + recall)

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

Results

Zero-shot prompting worked best overall

It caught all malware (100% recall) and only a few false alarms

Chain-of-thought and few-shot prompts also
did very well

Encouraging the model to "think step by step" helps

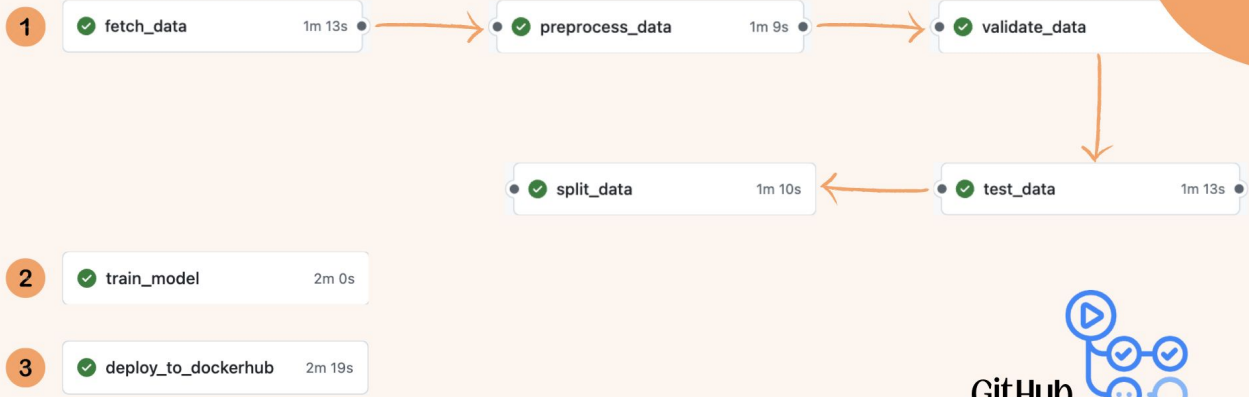
One-shot and negative prompts did worse

Missed many attacks but were cautious (few false alarms)

How can we use this?

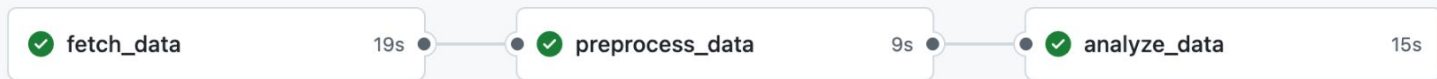
It can be an additional layer to our existing system

Automatization



data.yaml

on: push



Takeaways

- Even simple prompts ("**zero-shot**") worked very well here
- More guidance helps, but not always better than a simple question
- Cautious prompts lower false alarms but miss real attacks
- Bigger, more expensive models aren't always better: **GPT-4o** and **o1** gave same results, but **GPT-4o** was cheaper
- For best balance of accuracy and cost: Use **GPT-4o**
- For cheap bulk triage: Use **o4-mini** (accepting some misses)

Takeaways

- Don't fully rely on AI
- Be careful with sensitive data

Thank you!

Q&A